

Outline

- 1. Introduction**
2. Vectors and distances
3. Defining a new representation
4. Changing Coordinates
5. Dimension Reduction by compression
6. Conclusion, extensions
7. Alternatives to PCA, non linear embedding methods
8. Annexes
9. Principal Components and orthogonal subspaces

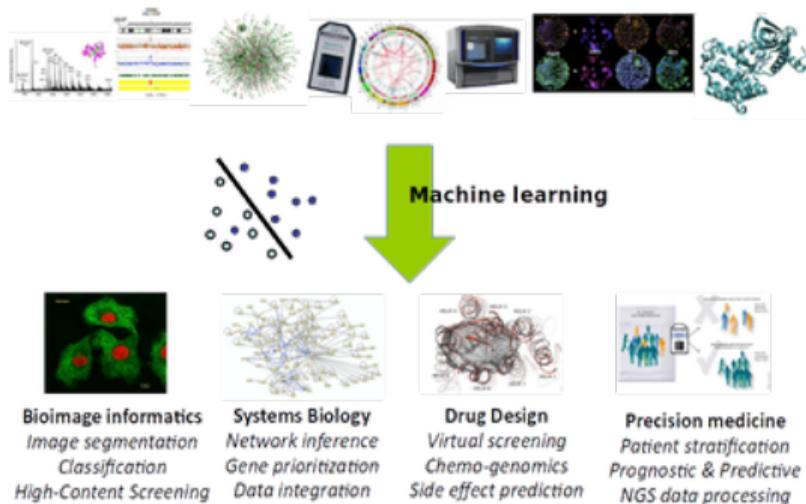
Are machines learning ?

- Convergence of math/info/computer science research
- Inspired by research in neuroscience and cognition (and science fiction)
- Contemporary to high throughput data acquisition
- Two basic ingredients: data and algorithms



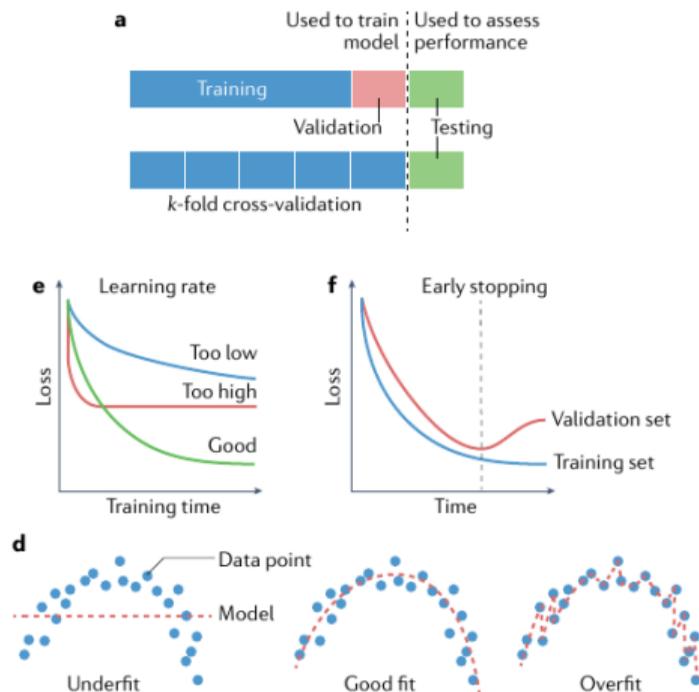
Machine learning in Biology

- Data have grown in complexity and size in all fields of biology
- Data are multimodal: sequences, structures, spectra, images, molecular, clinical, evolutionary
- Impossible to handle the analysis by descriptive methods only
- Reproducible research
- Machine Learning in Biology has become a field on its own



The quantitative shift [?]

- knowledge transfer of basic concepts in ML for biologists
- training / testing
- over-fitting / under-fitting
- Linear / non linear
- Interpretability of ML methods
- Computational complexity / time



The two sides of Machine Learning

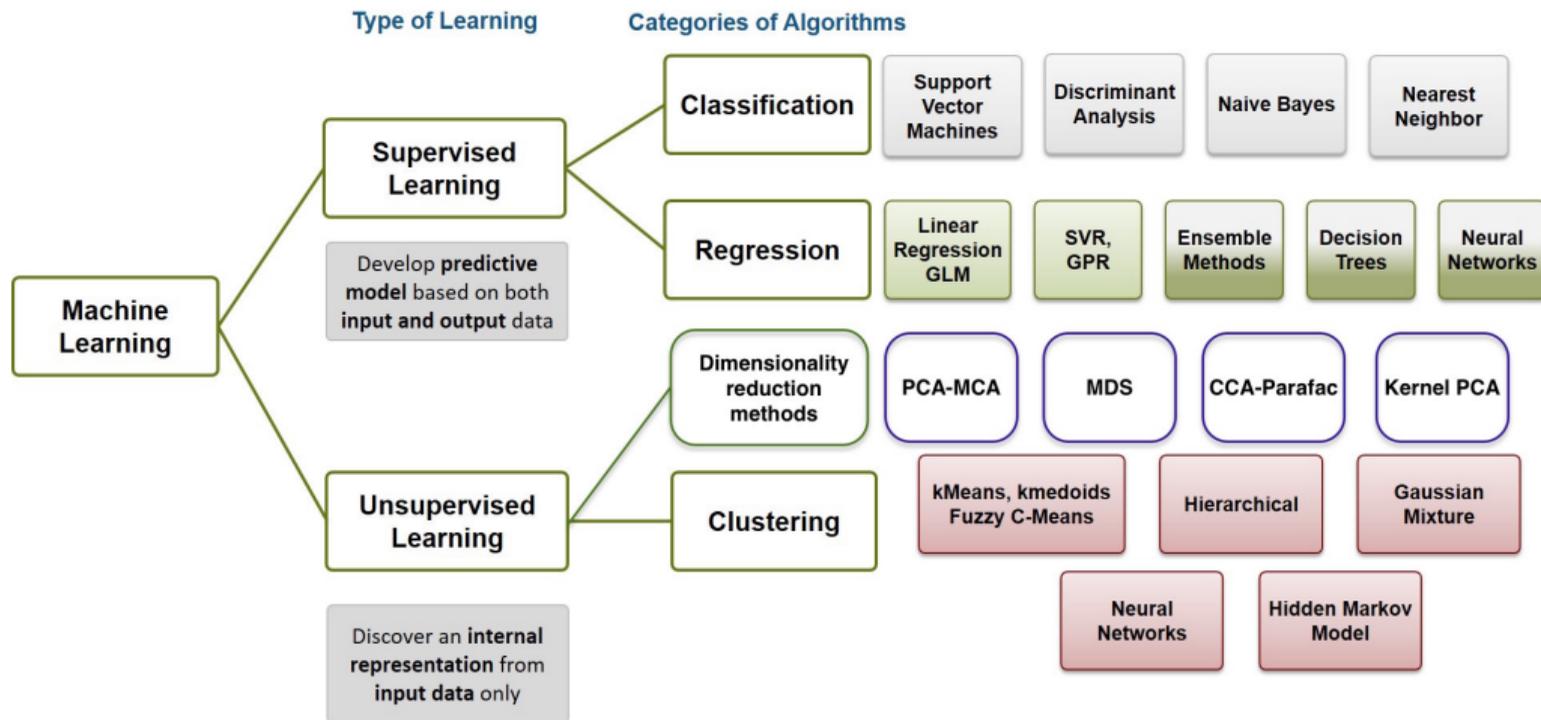
Supervised Learning

- Observe $(y_1, x_1), \dots, (y_n, x_n)$
 - Construct a predictor $f : \mathcal{X} \rightarrow \mathcal{Y}$
 - Define a loss function $\ell(y, f(x))$ to score predictions
 - Minimize the generalization error (new y ?)
- Regression, classification

Non-Supervised Learning

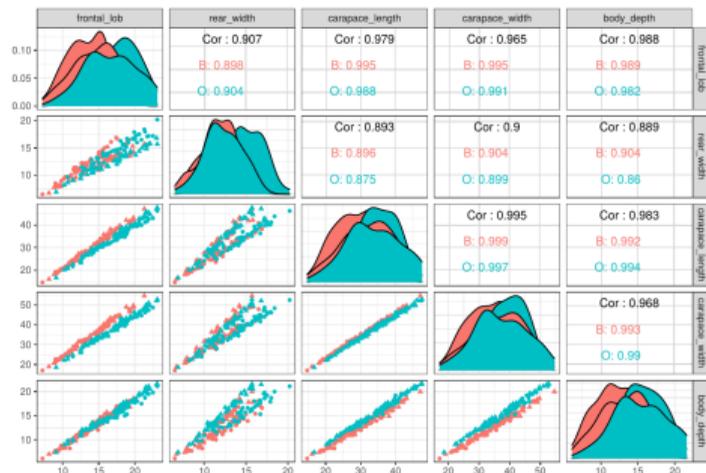
- Observe (x_1, \dots, x_n)
 - Describe the structure of X without external information
 - Group individuals ? Group variables ?
 - Loss is more difficult to define
- Dimension Reduction, Clustering

Rough typology of ML methods



The purpose of Dimension Reduction

- Visualization (> 2 variables)
- Multivariate analysis (beyond pairwise)
- Summary of the data
- Redundancy
- Reduce the data for downstream methods



Crabs dataset ($n = 200, p = 8$)

An unprecedented challenge

- Genomics was precursor for data representation and visualization
- Gene Expression data $\sim 30,000$ variables
- Recent single-cell technologies: up to 10^6 cells

Publication	cells	tissue	Seq. protocol
Cadwell et al. (2016)	46	visual cortex	Smart-seq2
Tasic et al. (2016)	1,679	visual cortex	SMARTer
Macosko et al. (2015)	44,808	retina	Drop-seq
10x Genomics	1,306,127	brain cells	10x Gen.Chrom.

- Dimension reduction is mandatory for any analysis (clustering, visualization, inference)

Outline

1. Introduction
- 2. Vectors and distances**
3. Defining a new representation
4. Changing Coordinates
5. Dimension Reduction by compression
6. Conclusion, extensions
7. Alternatives to PCA, non linear embedding methods
8. Annexes
9. Principal Components and orthogonal subspaces

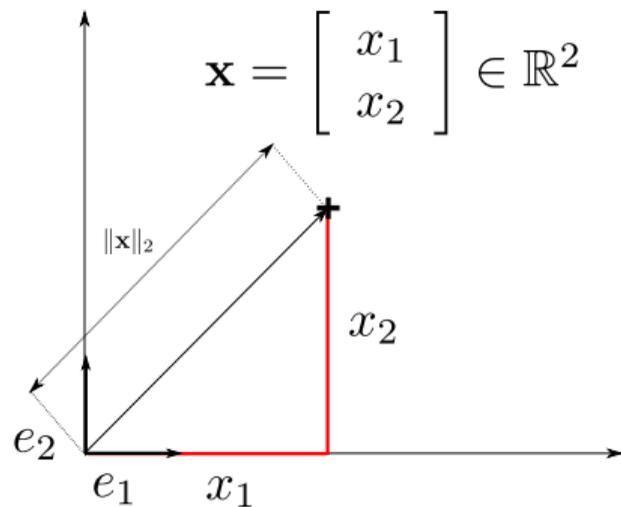
Vectors of \mathbb{R}^d

- \mathbf{x} is a vector of \mathbb{R}^d defined by a n -uplet (x_1, \dots, x_d) (coordinates)
- Considering the canonical basis ($d = 2$):

$$\mathbf{e}_1 = \begin{bmatrix} 1 \\ 0 \end{bmatrix}, \quad \mathbf{e}_2 = \begin{bmatrix} 0 \\ 1 \end{bmatrix}$$

- Its coordinates corresponds to a decomposition on the unitary basis:

$$\mathbf{x} \in \mathbb{R}^2, \quad \mathbf{x} = x_1 \begin{bmatrix} 1 \\ 0 \end{bmatrix} + x_2 \begin{bmatrix} 0 \\ 1 \end{bmatrix}$$

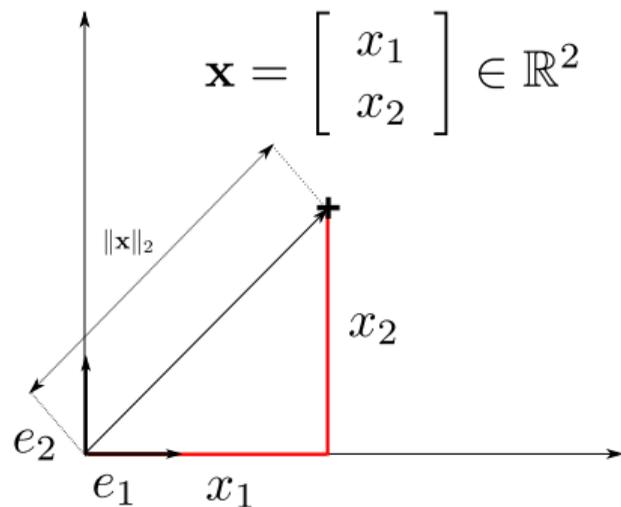


Vectors of \mathbb{R}^d

- A vector is a point in a space (here \mathbb{R}^2)
- Generalize for vectors of \mathbb{R}^d

$$\mathbf{x} \in \mathbb{R}^d, \quad \mathbf{x} = \sum_{h=1}^d x_h \mathbf{e}_h$$

- By default, \mathbf{x} is a column vector, \mathbf{x}' its transpose
- Concatenate p column vectors $[\mathbf{x}_1, \dots, \mathbf{x}_p]$.



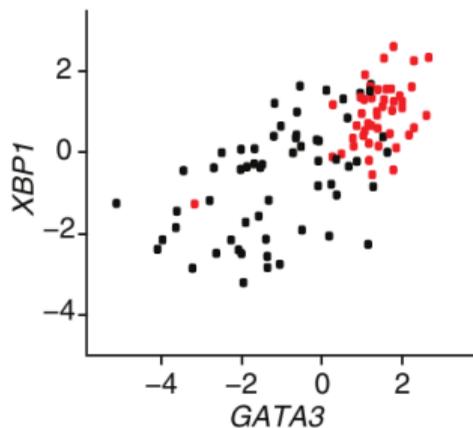
Expression of 105 breast tumor samples ER(+/-)

Measure of the expression of two genes XBP1 and GATA3 for $n = 105$ samples

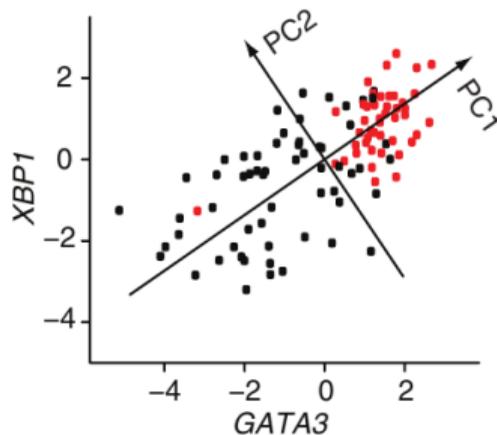
Each measure is denoted by $\mathbf{x}_i = \begin{bmatrix} x_i^{\text{GATA3}} \\ x_i^{\text{XBP1}} \end{bmatrix}$

Each point is a vector of \mathbb{R}^2

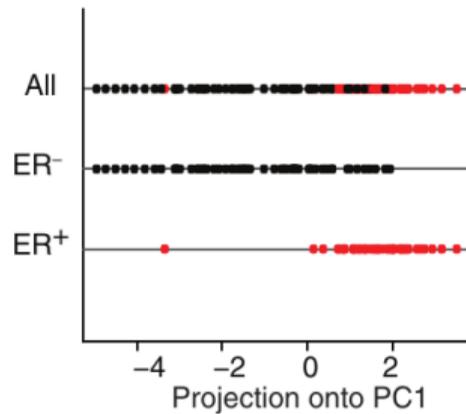
a



b



c



[?]

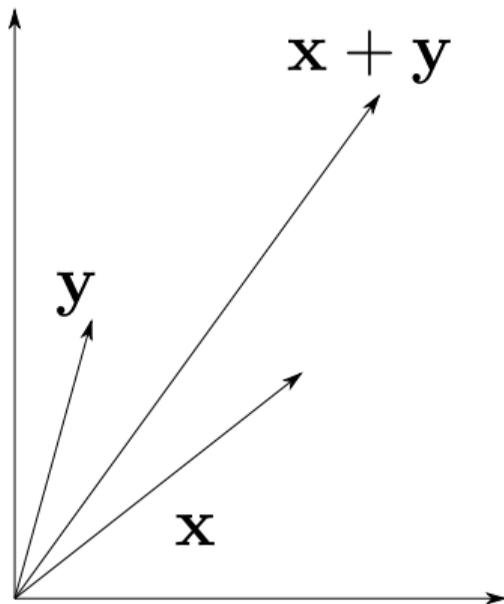
Vectors of \mathbb{R}^d and basic operations: Addition

$$\mathbf{x} + \mathbf{y} = \begin{bmatrix} x_1 + y_1 \\ \vdots \\ x_d + y_d \end{bmatrix}$$

$$(\mathbf{x} + \mathbf{y}) = (\mathbf{y} + \mathbf{x}),$$

$$(\mathbf{x} + \mathbf{y}) + \mathbf{z} = \mathbf{x} + (\mathbf{y} + \mathbf{z})$$

Associative, Commutative

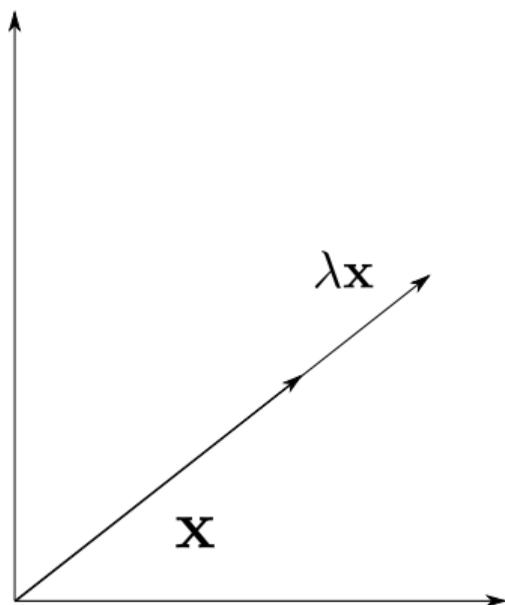


Vectors of \mathbb{R}^d and basic operations: Multiplication by a scalar

$$\forall \lambda \in \mathbb{R}, \lambda \mathbf{x} = \begin{bmatrix} \lambda x_1 \\ \vdots \\ \lambda x_d \end{bmatrix}$$
$$\lambda(\mathbf{x} + \mathbf{y}) = \lambda \mathbf{x} + \lambda \mathbf{y},$$

Linear Combination

$$(\lambda_1 + \lambda_2)\mathbf{x} = \lambda_1 \mathbf{x} + \lambda_2 \mathbf{x}$$



Dot Product between vectors

- The dot product \bullet between two vectors is defined by the sum of the products of all components

$$\mathbf{x} \bullet \mathbf{y} = \langle \mathbf{x}, \mathbf{y} \rangle = \mathbf{x}'\mathbf{y} = \sum_{i=1}^d x_i y_i, \quad \langle \mathbf{x}, \mathbf{y} \rangle = \langle \mathbf{y}, \mathbf{x} \rangle, \quad \langle \mathbf{x}, \mathbf{y} + \mathbf{z} \rangle = \langle \mathbf{x}, \mathbf{y} \rangle + \langle \mathbf{x}, \mathbf{z} \rangle$$

- The dot product between two vectors is a scalar
- Basic properties:

$$\langle \mathbf{x}, \mathbf{y} \rangle = \langle \mathbf{y}, \mathbf{x} \rangle, \quad \langle \mathbf{x}, \mathbf{y} + \mathbf{z} \rangle = \langle \mathbf{x}, \mathbf{y} \rangle + \langle \mathbf{x}, \mathbf{z} \rangle, \quad \lambda \langle \mathbf{x}, \mathbf{y} \rangle = \langle \lambda \mathbf{x}, \mathbf{y} \rangle = \langle \mathbf{x}, \lambda \mathbf{y} \rangle$$

Norm of a vector and basic properties

- The euclidean norm (length of a vector)

$$\|\mathbf{x}\|_2^2 = \langle \mathbf{x}, \mathbf{x} \rangle = \mathbf{x}'\mathbf{x} = \sum_{i=1}^d x_i^2$$

- Non negativity : $\|\mathbf{x}\|_2 \geq 0$
- Definiteness : $\|\mathbf{x}\|_2 = 0 \leftrightarrow \mathbf{x} = \mathbf{0}$
- Triangle Inequality : $\|\mathbf{x} + \mathbf{y}\|_2 \leq \|\mathbf{x}\|_2 + \|\mathbf{y}\|_2$
- Homogeneity : $\|\lambda \times \mathbf{x}\|_2 = |\lambda| \times \|\mathbf{x}\|_2, \quad \lambda \in \mathbb{R}$

Principal norms used in Machine Learning

- L^1 norm or Manhattan norm:

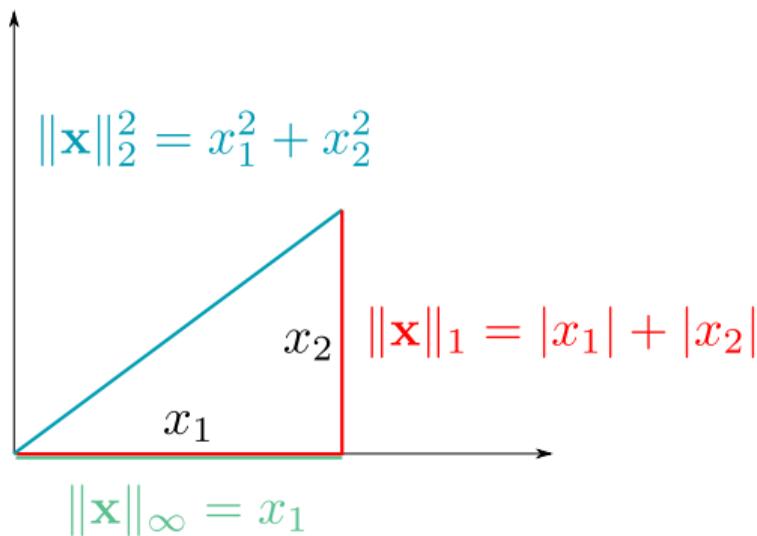
$$\|\mathbf{x}\|_1 = \sum_{i=1}^d |x_i|$$

- L^2 norm or Euclidian norm:

$$\|\mathbf{x}\|_2^2 = \sum_{i=1}^d x_i^2$$

- L^∞ norm or sup-norm:

$$\|\mathbf{x}\|_\infty = \max_{i=1, \dots, d} (|x_i|)$$



There are different ways to measure the norm of a vector

From norms to distances between vectors

- L^1 distance or Manhattan distance:

$$d_1(\mathbf{x}, \mathbf{y}) = \|\mathbf{x} - \mathbf{y}\|_1 = \sum_{i=1}^n |x_i - y_i|$$

- L^2 distance or Euclidean distance:

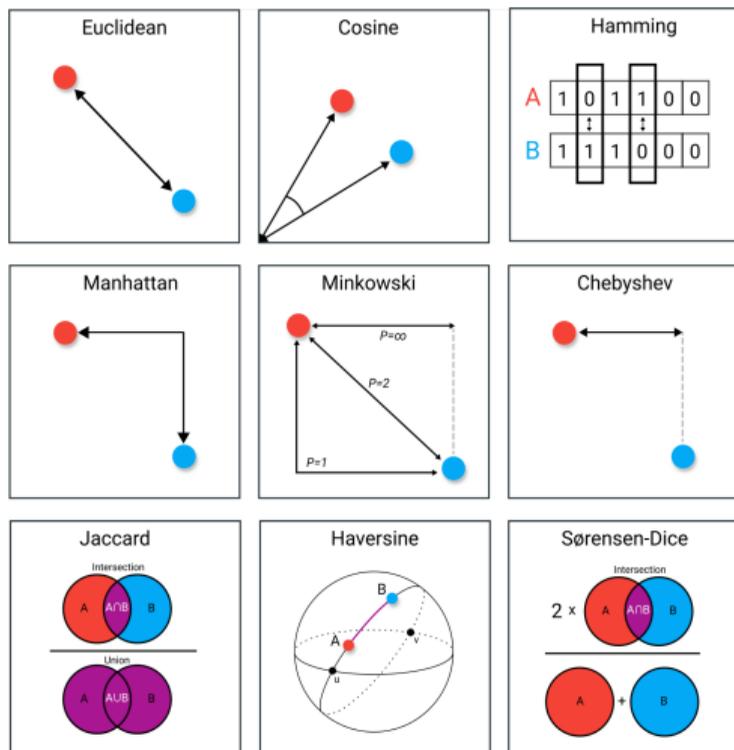
$$d_2^2(\mathbf{x}, \mathbf{y}) = \|\mathbf{x} - \mathbf{y}\|_2^2 = \sum_{i=1}^n (x_i - y_i)^2$$

- L^∞ distance or sup-distance:

$$d_\infty(\mathbf{x}, \mathbf{y}) = \|\mathbf{x} - \mathbf{y}\|_\infty = \max_{i=1, \dots, n} (|x_i - y_i|)$$

What drives the choice of a distance ?

- L^1 distance or Manhattan distance:
 - Adapted to discrete inputs
 - Robust to outliers
 - Non differentiable
- L^2 distance or Euclidean distance:
 - Most common, differentiable
 - Sensitive to dimension and outliers
 - Sensitive to the scale of the different inputs
- L^∞ distance or sup-distance:
 - Applied in logistical problems
 - More specific, less used



Dot product and orthogonal projection (1)

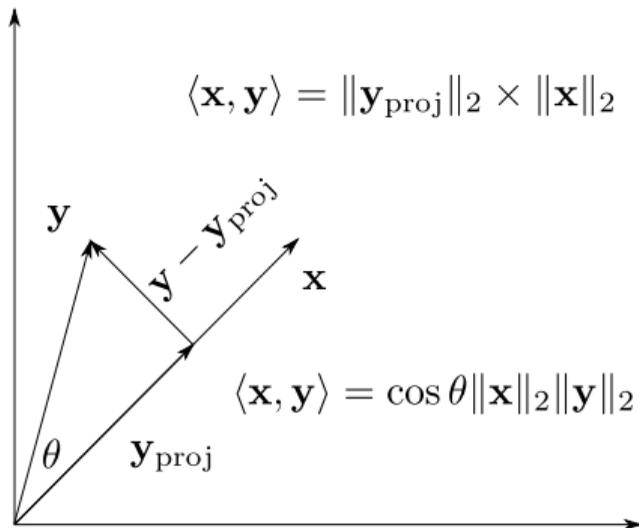
- The orthogonal projection of \mathbf{y} on \mathbf{x}

$$\mathbf{y}_{\text{proj}} = \lambda \mathbf{x} \quad , \quad \text{colinearity}$$

$$\mathbf{y} - \mathbf{y}_{\text{proj}} \perp \mathbf{x} \quad , \quad \text{orthogonality of residuals}$$

- The proportionality coefficient is given by

$$\lambda = \frac{\langle \mathbf{y}, \mathbf{x} \rangle}{\|\mathbf{x}\|_2}$$



$$\|\mathbf{y}\|_2^2 = \|\mathbf{y}_{\text{proj}}\|_2^2 + \|\mathbf{y} - \mathbf{y}_{\text{proj}}\|_2^2$$

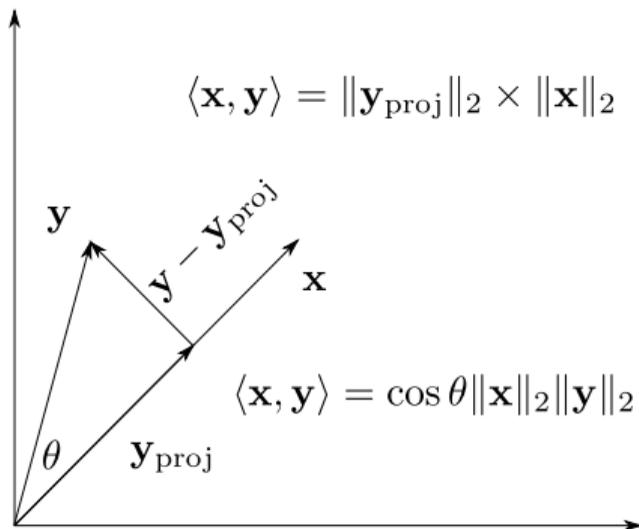
Dot product and orthogonal projection (2)

- Using trigonometry properties:

$$\cos \theta = \frac{\|\mathbf{y}_{proj}\|_2}{\|\mathbf{y}\|_2} = \lambda \frac{\|\mathbf{x}\|_2}{\|\mathbf{y}\|_2}$$

- The dot product is the length of \mathbf{x} times the length of the ortho. projection of \mathbf{y}
- Orthogonality :

$$\mathbf{x} \perp \mathbf{y} \leftrightarrow \langle \mathbf{y}, \mathbf{x} \rangle = 0$$



$$\|\mathbf{y}\|_2^2 = \|\mathbf{y}_{proj}\|_2^2 + \|\mathbf{y} - \mathbf{y}_{proj}\|_2^2$$

Outline

1. Introduction
2. Vectors and distances
- 3. Defining a new representation**
4. Changing Coordinates
5. Dimension Reduction by compression
6. Conclusion, extensions
7. Alternatives to PCA, non linear embedding methods
8. Annexes
9. Principal Components and orthogonal subspaces

Centering a dataset (1)

- The empirical mean of \mathbf{x}^j :

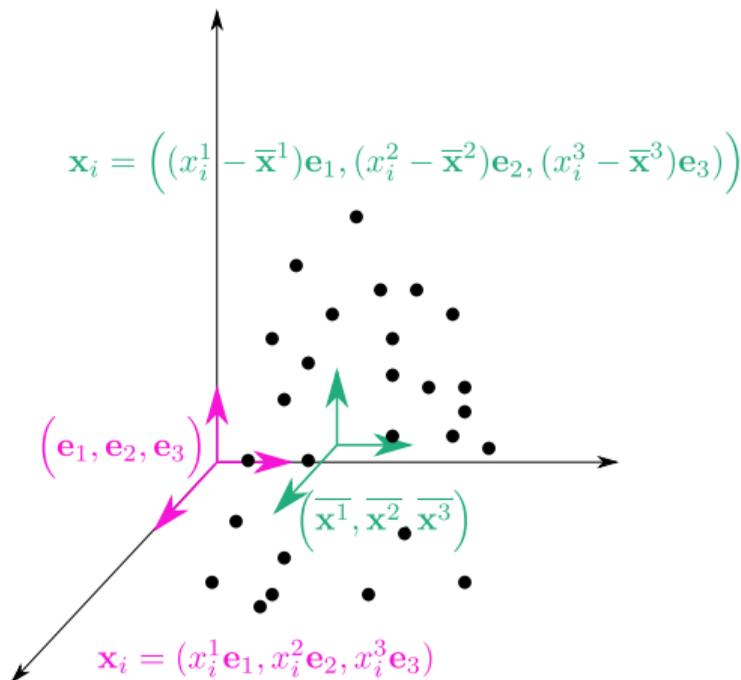
$$\bar{x}^j = \frac{1}{n} \sum_{i=1}^n x_i^j, \quad \bar{\mathbf{x}}_j = \bar{x}^j \times \mathbf{1}_n$$

- The empirical mean of \mathbf{x}^j is its projection on the constant

$$\bar{x}^j = \frac{1}{n} (1, \dots, 1) \bullet \mathbf{x} = \frac{1}{n} \langle \mathbf{1}'_n, \mathbf{x}^j \rangle$$

- The vector of means is the barycenter of the data

$$\bar{\mathbf{x}} = [\bar{x}^1, \dots, \bar{x}^p]$$

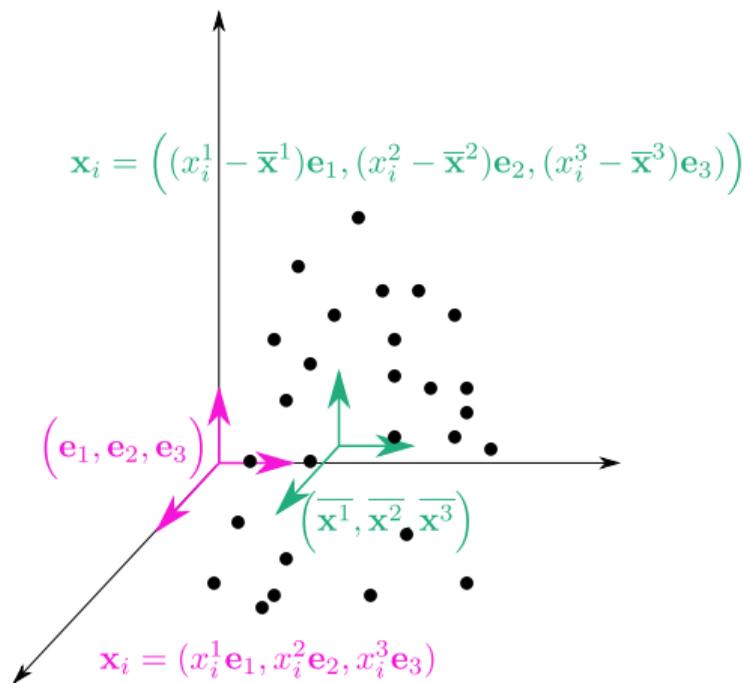


Centering a dataset: changing the origin

- Consists in removing the mean of each variable

$$\mathbf{X}_c = \left[\mathbf{x}^1 - \bar{\mathbf{x}}^1, \dots, \mathbf{x}^P - \bar{\mathbf{x}}^P \right]$$

- Centering to avoid positional effects
- $\bar{\mathbf{x}}$ becomes the new origin



Scaling a dataset

- The empirical variance of \mathbf{x}^j :

$$\text{var}(\mathbf{x}^j) = \frac{1}{n} \sum_{i=1}^n (x_i^j - \bar{x}^j)^2$$

- It is the distance of variable \mathbf{x}_j to its mean

$$\text{var}(\mathbf{x}^j) = \frac{1}{n} \|\mathbf{x}^j - \bar{\mathbf{x}}^j\|_2^2 = \frac{1}{n} \langle \mathbf{x}^j - \bar{\mathbf{x}}^j, \mathbf{x}^j - \bar{\mathbf{x}}^j \rangle = \frac{1}{n} \mathbf{x}_c^j \bullet \mathbf{x}_c^j$$

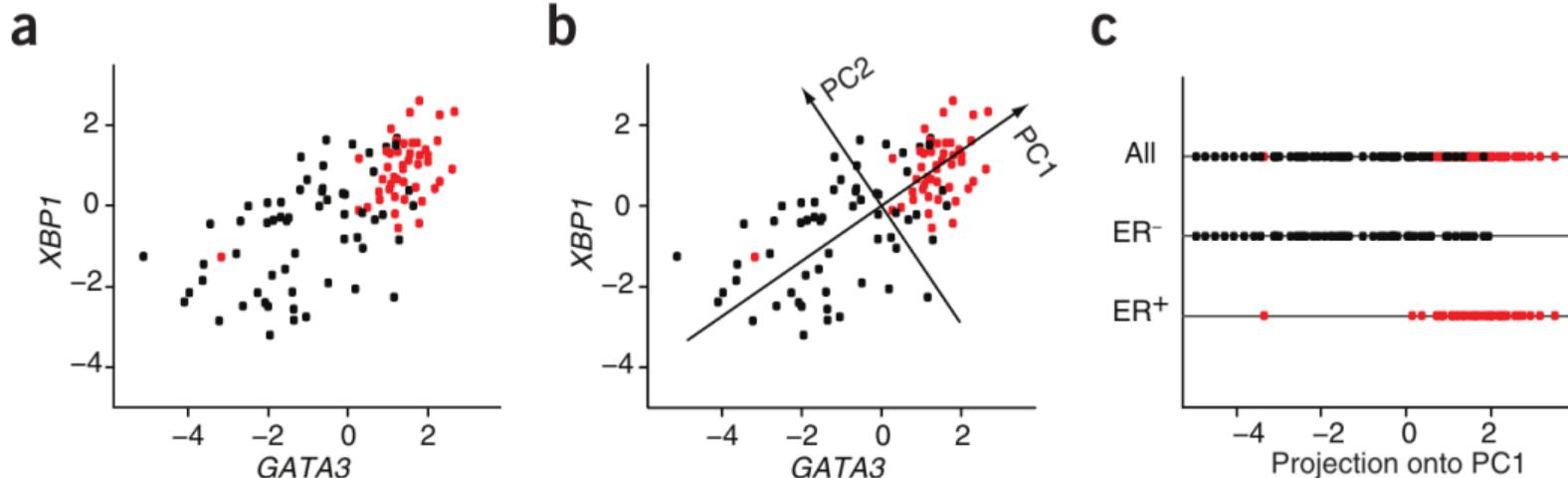
- The empirical variance is the length of the residuals (after centering)
- Scaling to standardize variables contributions (unitary variance)

$$\tilde{\mathbf{X}}_c = \left[\frac{\mathbf{x}^1 - \bar{\mathbf{x}}^1}{\text{var}^{1/2}(\mathbf{x}^1)}, \dots, \frac{\mathbf{x}^p - \bar{\mathbf{x}}^p}{\text{var}^{1/2}(\mathbf{x}^p)} \right]$$

Expression of 105 breast tumor samples ER(+/-)

The data matrix is

$$\mathbf{x}_c = \left[\mathbf{x}_c^{\text{GATA3}}, \mathbf{x}_c^{\text{XBP1}} \right]_{105 \times 2}$$



The expression of those 2 genes is very correlated: redundancy between columns

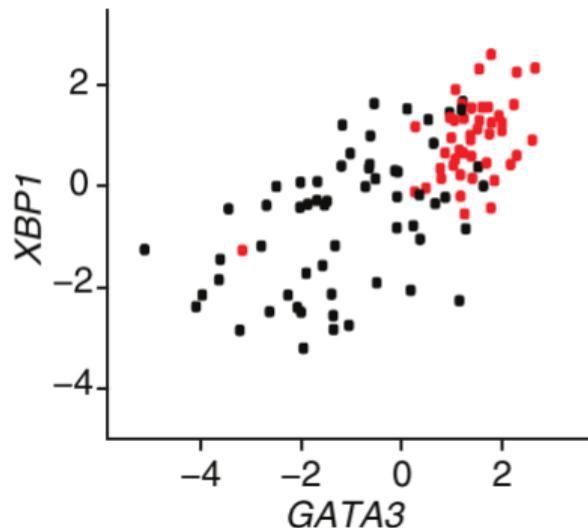
Covariance and Correlation between variables

- The empirical covariance between variables

$$c(\mathbf{x}^j, \mathbf{x}^{j'}) = \frac{1}{n} \sum_{i=1}^n (x_i^j - \bar{x}^j)(x_i^{j'} - \bar{x}^{j'})$$

$$r(\mathbf{x}^j, \mathbf{x}^{j'}) = \frac{c(\mathbf{x}^j, \mathbf{x}^{j'})}{\sqrt{\text{var}(\mathbf{x}^j) \text{var}(\mathbf{x}^{j'})}}$$

- Quantifies the expected co-variations between variables
- If $r(\mathbf{x}^j, \mathbf{x}^{j'}) \simeq 1$ the two variables provide the same information



Distance and covariance

- Between-variables distance:

$$\begin{aligned}\frac{1}{n} \|\mathbf{x}_c^j - \mathbf{x}_c^{j'}\|^2 &= \frac{1}{n} \|\mathbf{x}_c^j\|^2 + \frac{1}{n} \|\mathbf{x}_c^{j'}\|^2 - 2 \frac{1}{n} \langle \mathbf{x}_c^j, \mathbf{x}_c^{j'} \rangle \\ &= \text{var}(\mathbf{x}^j) + \text{var}(\mathbf{x}^{j'}) - 2c(\mathbf{x}^{j'}, \mathbf{x}^j)\end{aligned}$$

- Normalized distance using centered and scaled variables

$$\frac{1}{n} \|\tilde{\mathbf{x}}_c^j - \tilde{\mathbf{x}}_c^{j'}\|^2 = 2 - 2r(\mathbf{x}^{j'}, \mathbf{x}^j)$$

- The correlation coefficient is a distance measure between variables:

$$r(\mathbf{x}^{j'}, \mathbf{x}^j) = 1 - \frac{1}{2} \times \frac{1}{n} \|\tilde{\mathbf{x}}_c^j - \tilde{\mathbf{x}}_c^{j'}\|^2$$

Correlation and distance between variables

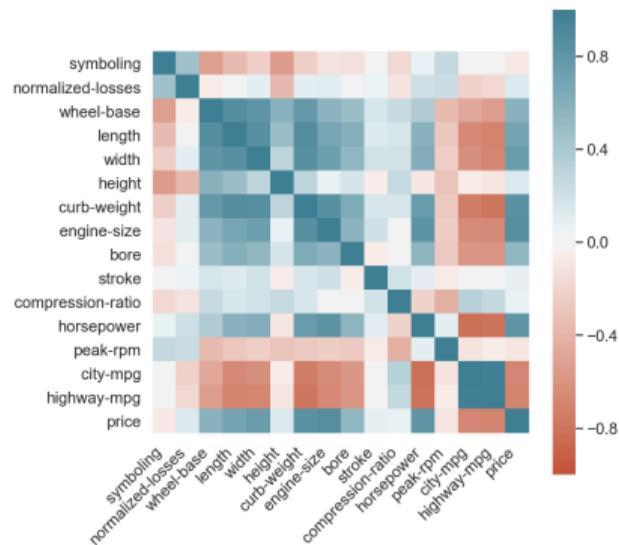
- Pairwise distance between variables

$$\mathbf{S} = \begin{bmatrix} c(\mathbf{x}^1, \mathbf{x}^1) & \dots & c(\mathbf{x}^{j'}, \mathbf{x}^j) \\ & \ddots & \\ c(\mathbf{x}^j, \mathbf{x}^{j'}) & \dots & c(\mathbf{x}^P, \mathbf{x}^P) \end{bmatrix}$$

- Normalized distance: correlation matrix

$$\mathbf{R} = \begin{bmatrix} r(\mathbf{x}^1, \mathbf{x}^1) & \dots & r(\mathbf{x}^{j'}, \mathbf{x}^j) \\ & \ddots & \\ r(\mathbf{x}^j, \mathbf{x}^{j'}) & \dots & r(\mathbf{x}^P, \mathbf{x}^P) \end{bmatrix}$$

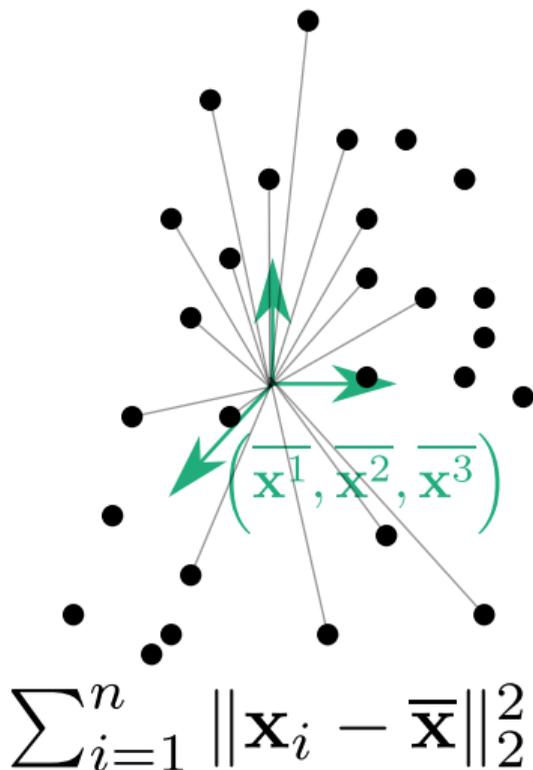
- Symmetric, invertible (semi definite positive)



Total Inertia of a dataset

The global variance of a dataset for centered variables

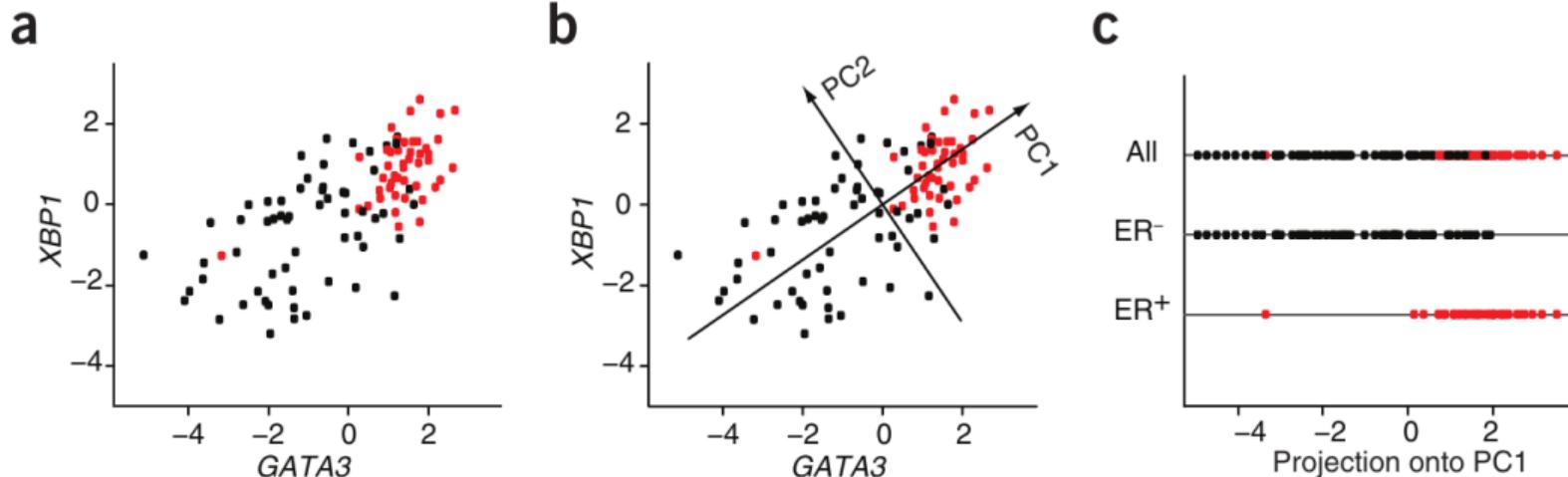
$$\begin{aligned} I_T(\mathbf{X}) &= \frac{1}{n} \sum_{i=1}^n \sum_{j=1}^p (x_i^j - \bar{x}^j)^2 \\ &= \sum_{j=1}^p \text{var}(\mathbf{x}^j) \end{aligned}$$



Inertia of a dataset

To generalize the notion of dispersion to a complete dataset:

$$I_T(\mathbf{X}) = \frac{1}{n} \sum_{i=1}^n \sum_{j=1}^p (x_i^j - \bar{x}^j)^2$$

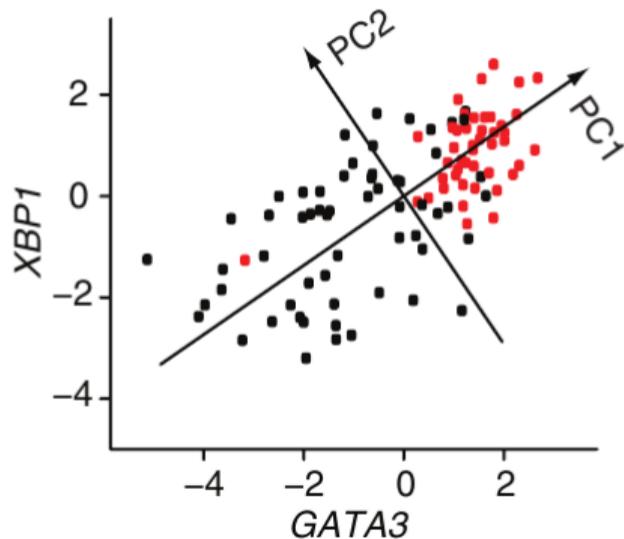


Outline

1. Introduction
2. Vectors and distances
3. Defining a new representation
- 4. Changing Coordinates**
5. Dimension Reduction by compression
6. Conclusion, extensions
7. Alternatives to PCA, non linear embedding methods
8. Annexes
9. Principal Components and orthogonal subspaces

Outline before PCA

- PCA is based on a change in coordinates
- Before performing PCA, focus on the rotation of a dataset
- Change coordinates from 2D to 2D, then generalize



From 2D to 2D with rotation

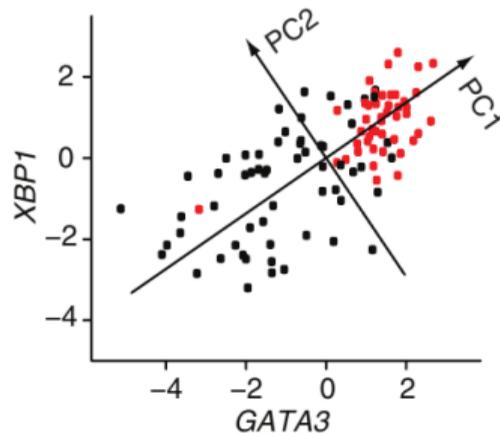
- Find new coordinates \mathbf{Z} to better represent \mathbf{X}
- Define z_{1i} the new coordinates of individual i on axis 1 as linear combinations of the ancient coordinates

$$z_{1i} = v_{11}\tilde{x}_{i,c}^1 + v_{12}\tilde{x}_{i,c}^2$$

- This operation resumes to a linear transform of \mathbf{x}_i (old) to obtain \mathbf{z} (new)

$$\mathbf{z}_{i1} = \tilde{\mathbf{x}}_{i,c} \mathbf{v}_1$$

- How to determine $\mathbf{v}_1 = \begin{bmatrix} v_{11} \\ v_{12} \end{bmatrix}_{2 \times 1}$?



New Coordinates

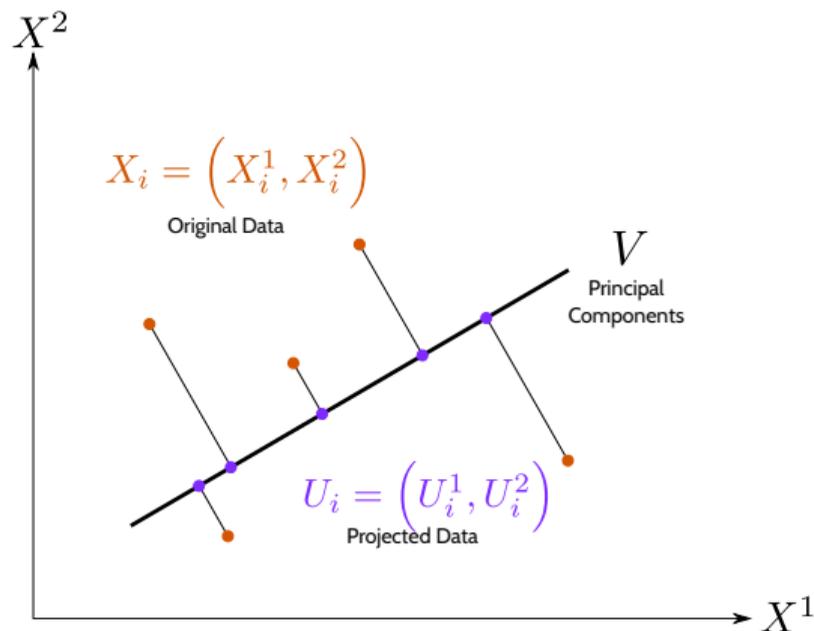
- In the example:

$$z_i^1 = 0.83 \times \text{GATA3}_i + 0.56 \times \text{XBP1}_i$$

- For the best representation of \mathbf{X}

$$\hat{v}_{11} = 0.83, \quad \hat{v}_{12} = 0.56,$$

- Notation \hat{v} stands for optimized coordinates



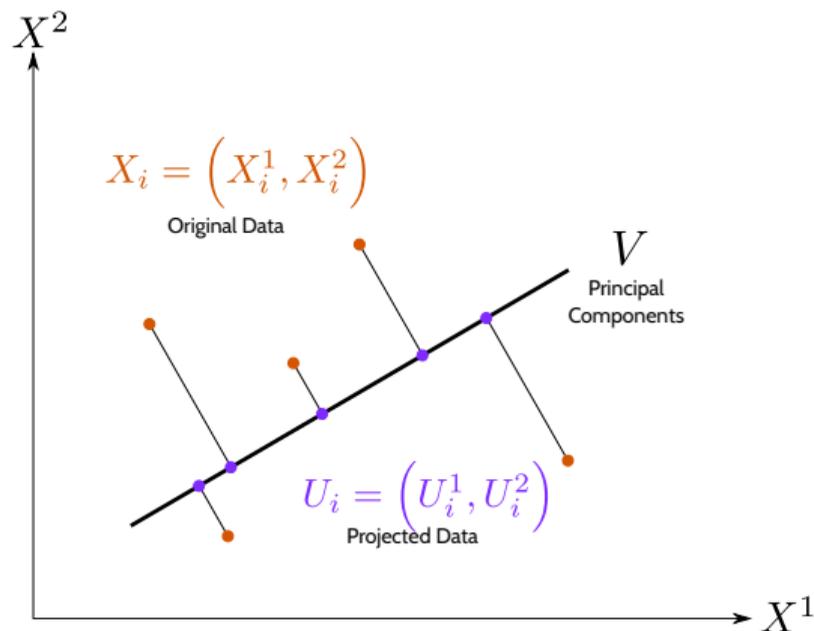
New coordinates in the matricial framework

- The coefficients are common to all individuals:

$$\begin{aligned} \mathbf{z}_1 &= v_{11}\tilde{\mathbf{x}}_c^1 + v_{12}\tilde{\mathbf{x}}_c^2 \\ &= \begin{bmatrix} \tilde{\mathbf{x}}_c^1 & \tilde{\mathbf{x}}_c^2 \end{bmatrix}_{n \times 2} \begin{bmatrix} v_{11} \\ v_{12} \end{bmatrix}_{2 \times 1} \end{aligned}$$

$$\mathbf{z}_1 = \tilde{\mathbf{X}}_c \mathbf{v}_1$$

- Equation of a line with slope \mathbf{v}_1
- Centered data so no intercept



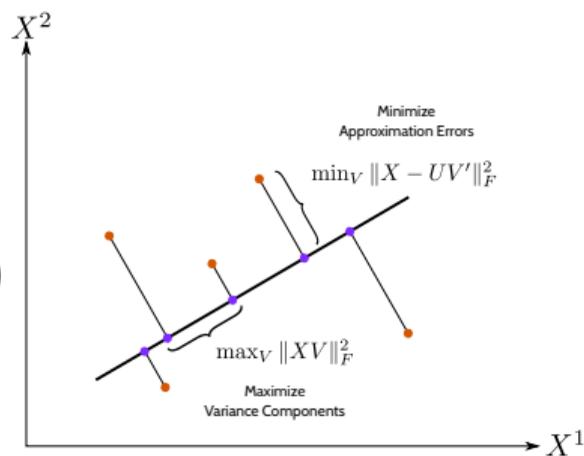
New coordinates in the matrix framework (1)

- First axis carries the biggest empirical variance

$$\begin{aligned}\text{var}(\mathbf{z}_1) &= \text{var}(\tilde{\mathbf{X}}_c \mathbf{v}_1) \\ &= \text{var}(v_{11} \tilde{\mathbf{x}}_c^1 + v_{12} \tilde{\mathbf{x}}_c^2) \\ &= v_{11}^2 \text{var}(\tilde{\mathbf{x}}_c^1) + v_{12}^2 \text{var}(\tilde{\mathbf{x}}_c^2) + 2v_{11}v_{12} c(\tilde{\mathbf{x}}_c^1, \tilde{\mathbf{x}}_c^2)\end{aligned}$$

- Using the standardized version (scaled)

$$\text{var}(\mathbf{z}_1) = v_{11}^2 + v_{12}^2 + 2v_{11}v_{12} \times r(\tilde{\mathbf{x}}_c^1, \tilde{\mathbf{x}}_c^2)$$

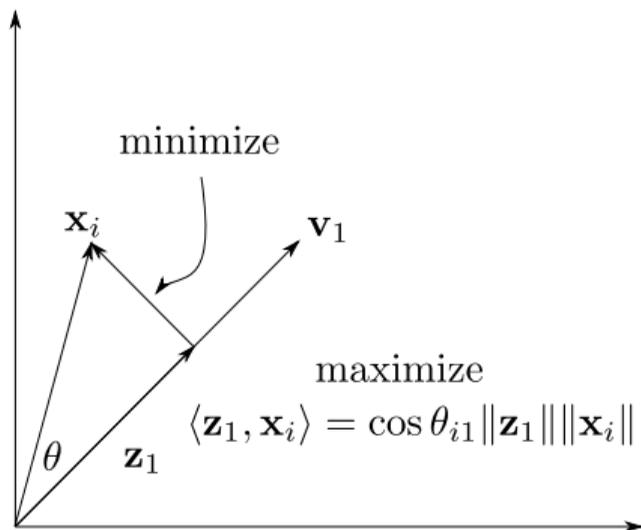


New coordinates in the matricial framework (2)

- To find the new coordinates: find \mathbf{v}_1 such that $\text{var}(\mathbf{z}_1)$ is maximal

$$\text{var}(\mathbf{z}_1) = v_{11}^2 + v_{12}^2 + 2v_{11}v_{12} \times r(\tilde{\mathbf{x}}_c^1, \tilde{\mathbf{x}}_c^2)$$

- Constraint for a normed basis: $\|\mathbf{v}_1\|_2^2 = 1$
- This ensures that the new basis is of unitary scale, so that the information carried by the new axes can be compared



PCA as an optimization problem

- To find the first axis, find coefficients \mathbf{v}_1 , s.t.

$$\begin{aligned}\max_{\mathbf{v}_1, \|\mathbf{v}_1\|_2^2=1} \left\{ \text{var}(\mathbf{z}_1) \right\} &= \max_{\mathbf{v}_1, \|\mathbf{v}_1\|_2^2=1} \left\{ \text{var}(\mathbf{X}_c \mathbf{v}_1) \right\} \\ &= \max_{\mathbf{v}_1, \|\mathbf{v}_1\|_2^2=1} \left\{ \mathbf{v}_1' (\mathbf{X}_c' \mathbf{X}_c) \mathbf{v}_1 \right\} \\ &= \max_{\mathbf{v}_1, \|\mathbf{v}_1\|_2^2=1} \left\{ \mathbf{v}_1' \mathbf{S} \mathbf{v}_1 \right\}\end{aligned}$$

- The solution of this optimization problem is explicit

$$\begin{aligned}\mathbf{v}_1' \mathbf{v}_1 &= 1 \\ \mathbf{S} \mathbf{v}_1 &= \lambda_1 \mathbf{v}_1\end{aligned}$$

- \mathbf{v}_1 (resp λ_1) is the first eigenvector (resp eigenvalue) of the **covariance** matrix

normed PCA as an optimization problem

- To find the first axis, find coefficients $\tilde{\mathbf{v}}_1$, s.t.

$$\begin{aligned}\max_{\tilde{\mathbf{v}}_1, \|\tilde{\mathbf{v}}_1\|_2=1} \left\{ \text{var}(\mathbf{z}_1) \right\} &= \max_{\tilde{\mathbf{v}}_1, \|\tilde{\mathbf{v}}_1\|_2=1} \left\{ \text{var}(\tilde{\mathbf{X}}_c \tilde{\mathbf{v}}_1) \right\} \\ &= \max_{\tilde{\mathbf{v}}_1, \|\tilde{\mathbf{v}}_1\|_2=1} \left\{ \tilde{\mathbf{v}}_1 \left(\tilde{\mathbf{X}}_c' \tilde{\mathbf{X}}_c \right) \tilde{\mathbf{v}}_1 \right\} \\ &= \max_{\tilde{\mathbf{v}}_1, \|\tilde{\mathbf{v}}_1\|_2=1} \left\{ \tilde{\mathbf{v}}_1 \mathbf{R} \tilde{\mathbf{v}}_1' \right\}\end{aligned}$$

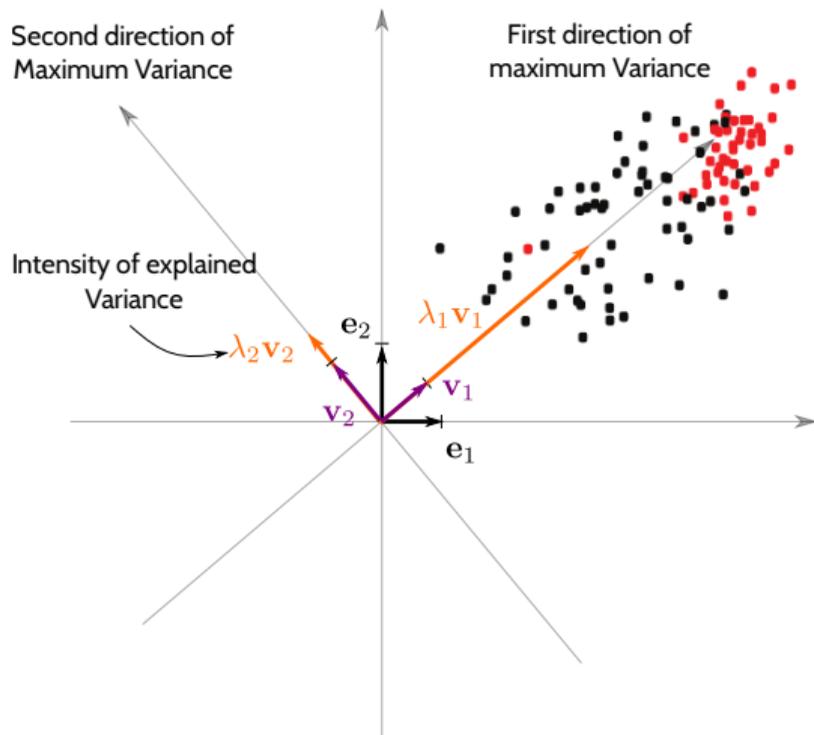
- The solution of this optimization problem is explicit

$$\begin{aligned}\tilde{\mathbf{v}}_1' \tilde{\mathbf{v}}_1 &= 1 \\ \mathbf{R} \tilde{\mathbf{v}}_1 &= \lambda_1 \tilde{\mathbf{v}}_1\end{aligned}$$

- $\tilde{\mathbf{v}}_1$ (resp λ_1) is the first eigenvector (resp eigenvalue) of the **correlation** matrix

Eigen Representation of the data

- \mathbf{S} contains the directions of maximal variance of the data
- $\mathbf{v}_1 \perp \mathbf{v}_2$ and are normed (unit variance)
- (λ_1, λ_2) quantify the amount of variance in each direction
- The eigen decomposition provides the best representation of the data in terms of variance
- Its the linear transform that makes the new set of coordinates diagonal



Quality of the representation

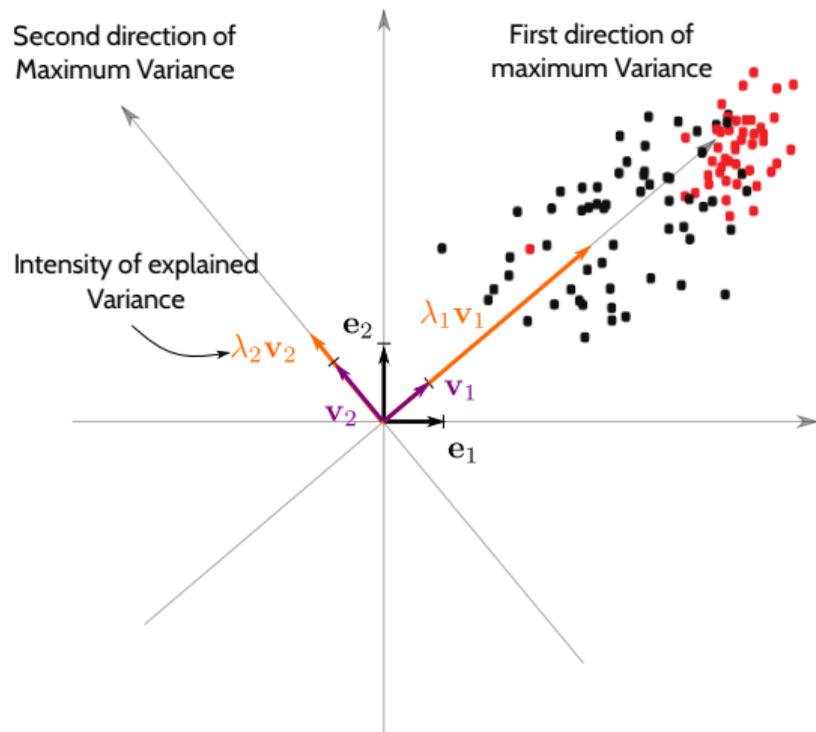
- Eigenvalues quantify the inertia of the dataset:

$$I_T(X) = \sum_{k=1} I_k(X) = \sum_{k=1}^K \lambda_k$$

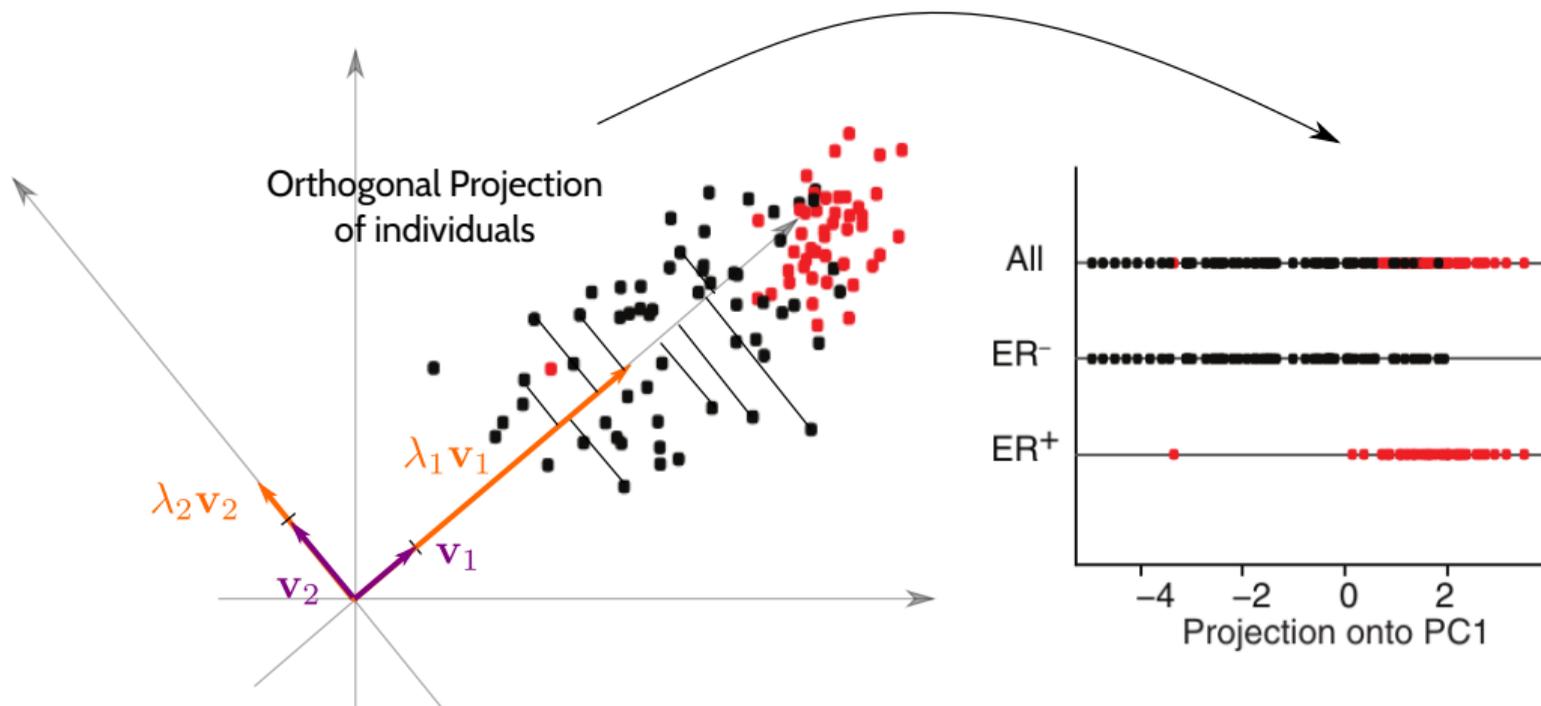
- Percent of explained variance:

$$\text{Contrib}_k = \frac{\lambda_k}{\sum_{\ell=1}^K \lambda_\ell}$$

$$\text{Contrib}_{1:k} = \frac{\sum_{h=1}^k \lambda_h}{\sum_{\ell=1}^K \lambda_\ell}$$



Representation of individuals in the new coordinates



The new coordinates for individuals are $(\mathbf{x}_i - \bar{\mathbf{x}}) \mathbf{v}_k$

Outline

1. Introduction
2. Vectors and distances
3. Defining a new representation
4. Changing Coordinates
- 5. Dimension Reduction by compression**
6. Conclusion, extensions
7. Alternatives to PCA, non linear embedding methods
8. Annexes
9. Principal Components and orthogonal subspaces

Outline

- In a first step we changed coordinates for better representation
- From 2D to 2D, there is no dimension reduction !
- The approach is generalized from p variables to K principal components

$$\mathbf{z}_k = \sum_{j=1}^p v_{kj} \tilde{\mathbf{x}}_c^j = \mathbf{X}_c \mathbf{v}_1$$

- Intuition: if v_{kj} is high, variable j highly contributes to principal component \mathbf{z}_k
- From p to $K(= 2)$ the information was compressed

General Case with K principal components

- $\mathbf{V}_{[p \times K]} = [\mathbf{v}_1, \dots, \mathbf{v}_K]$, the eigen vectors of the covariance matrix

$$\mathbf{S}_{p \times p} = \frac{1}{n} \mathbf{X}' \mathbf{X} = \frac{1}{n} \sum_{k=1}^K \lambda_k \mathbf{v}_k \mathbf{v}_k'$$

- $\mathbf{U}_{[n \times K]} = [\mathbf{u}_1, \dots, \mathbf{u}_K]$, the eigen vectors of the Gram matrix

$$\mathbf{G}_{n \times n} = \frac{1}{p} \mathbf{X} \mathbf{X}' = \frac{1}{p} \sum_{k=1}^K \lambda_k \mathbf{u}_k \mathbf{u}_k'$$

- Then we have

$$(\mathbf{X} \mathbf{X}') \mathbf{u}_k = \sqrt{\lambda_k} \mathbf{X} \mathbf{v}_k = \lambda_k \mathbf{u}_k$$

$$(\mathbf{X}' \mathbf{X}) \mathbf{v}_k = \sqrt{\lambda_k} \mathbf{X}' \mathbf{u}_k = \lambda_k \mathbf{v}_k$$

Low-rank approximation of \mathbf{X}

- The rank of a matrix (r^*) is the number of linearly independent columns (unknown in practice)
- From a statistical perspective, it is the number of independent coordinates that can describe a dataset
- The initial dataset can be rewritten such that

$$\mathbf{X} = \mathbf{U}_{n \times r^*} \mathbf{V}'_{r^* \times p} = \sum_{k=1}^{r^*} \sqrt{\lambda_k} \mathbf{u}_k \mathbf{v}'_k$$

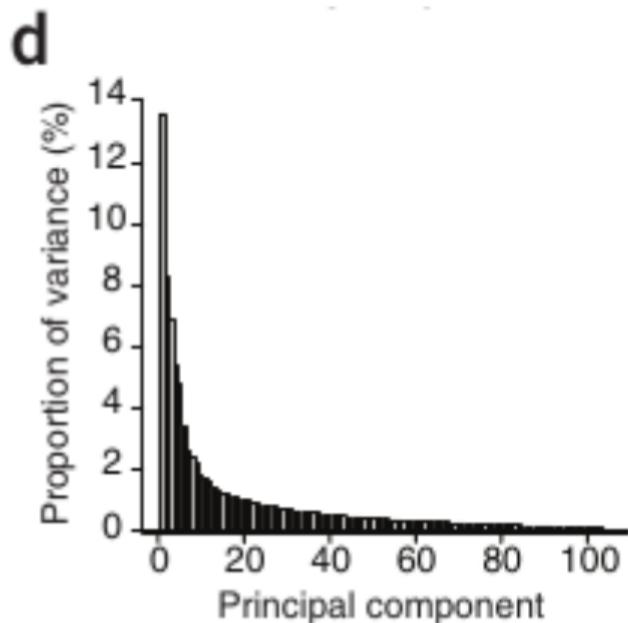
- Since the rank is unknown, we select a number of components K , and then:

$$\mathbf{X} \simeq \mathbf{U}_{n \times K} \mathbf{V}'_{K \times p} = \sum_{k=1}^K \sqrt{\lambda_k} \mathbf{u}_k \mathbf{v}'_k$$

- It is called the low-rank approximation of \mathbf{X}

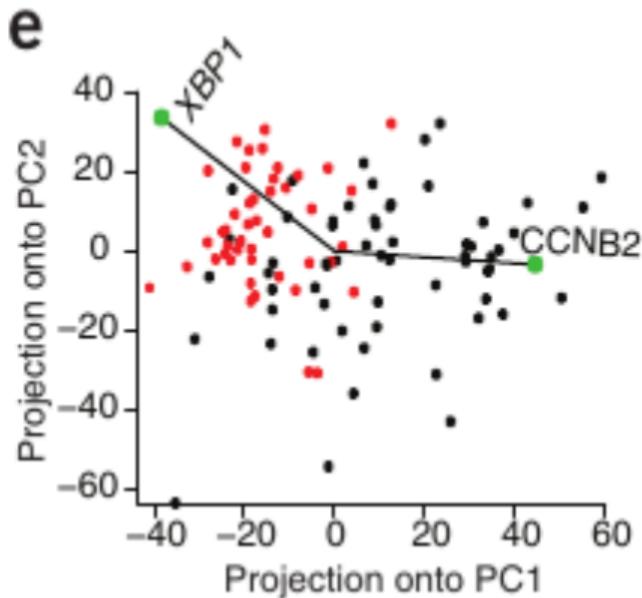
PCA on the complete ER dataset - 1

- First examples on 2 genes without dimension reduction
- PCA on the $p = 8534$ genes, $n = 105$ individuals
- $K_{\max} = 8534$ possible eigenvectors
- $\text{Contrib}_{1:2} \simeq 22\%$
- $\text{Contrib}_{1:63} \simeq 90\%$
- $\text{Contrib}_{1:104} \simeq 100\%$
- Choosing 104 eigenvectors reduces the dimension without too much loss
- Dimension reduction : from 8534 original variables to 104 new variables



PCA on the complete ER dataset - 2

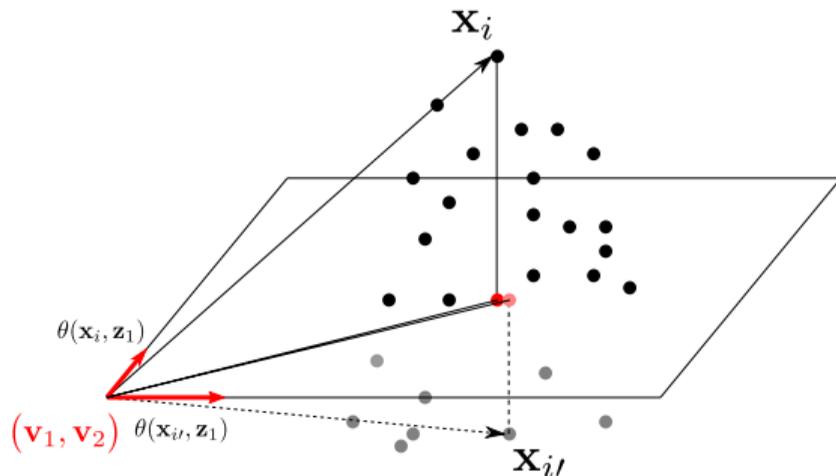
- Represent the data in the new coordinates (PCs)
- In the example the clusters (ER+/ER-) are more separable in the new representation
- Identify the contribution of genes to the axes
- Essential to interpret the new representation



Quality of the representation of individuals

- An individual \mathbf{x}_i is well represented if it is close to the axis \mathbf{z}_k
- Geometrically, $\mathbf{x}_i - \bar{\mathbf{x}}$ is colinear to \mathbf{z}_k
- Compute

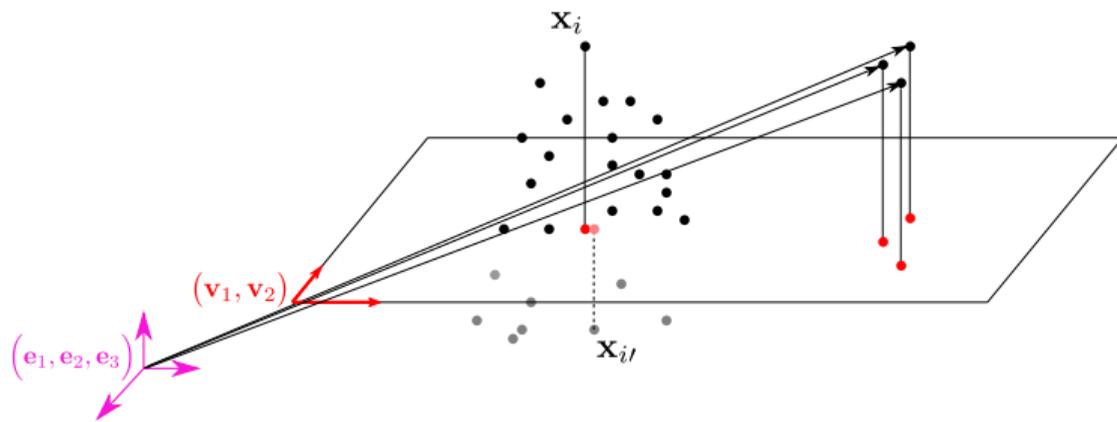
$$\cos^2 \theta(\mathbf{x}_i - \bar{\mathbf{x}}, \mathbf{z}_k) = \frac{\left((\mathbf{x}_i - \bar{\mathbf{x}}) \mathbf{v}_k \right)^2}{\|\mathbf{x}_i - \bar{\mathbf{x}}\|^2 \|\mathbf{v}_k\|^2}$$



Contribution of individuals to the representation

The contribution of a \mathbf{x}_i is the proportion of carried by \mathbf{x}_i

$$\text{contr}(\mathbf{x}_i, \mathbf{z}_k) = \frac{\left((\mathbf{x}_i - \bar{\mathbf{x}}) \mathbf{v}_k \right)^2}{n \lambda_k}$$



Properties of Principal components: the variable point of view

- Start with p correlated (redundant) variables $\tilde{\mathbf{X}}_c = \left[\tilde{\mathbf{x}}_c^1, \dots, \tilde{\mathbf{x}}_c^p \right]$ with

$$\mathbf{R}_{p \times p} = \begin{bmatrix} r(\mathbf{x}^1, \mathbf{x}^1) & \dots & r(\mathbf{x}^{j'}, \mathbf{x}^j) \\ & \ddots & \\ r(\mathbf{x}^j, \mathbf{x}^{j'}) & \dots & r(\mathbf{x}^p, \mathbf{x}^p) \end{bmatrix} = \frac{1}{n} \tilde{\mathbf{X}}_c' \tilde{\mathbf{X}}_c = \sum_{k=1}^K \lambda_k \mathbf{v}_k \mathbf{v}_k'$$

- Get K new uncorrelated (non redundant) variables $\mathbf{Z} = \left[\mathbf{z}^1, \dots, \mathbf{z}^K \right]$

Correlation Circle

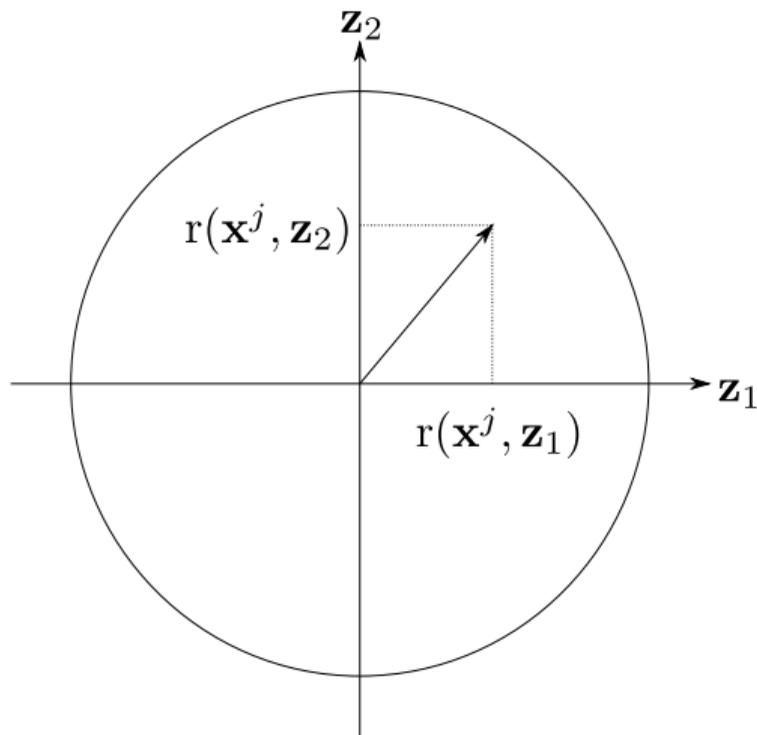
- Components are independent of variance with $\text{var}(\mathbf{z}_k) = \lambda_k$

$$\mathbf{S}_Z = \begin{bmatrix} \lambda_1 & & 0 \\ & \ddots & \\ 0 & & \lambda_K \end{bmatrix}$$

- Contribution of variables to axis:

$$\begin{aligned} c(\mathbf{x}^j, \mathbf{z}_k) &= (\mathbf{x}^j)' \mathbf{u}_k = \lambda_k v_{jk} \\ &= r(\mathbf{x}^j, \mathbf{z}_k) \text{ for normed PCA} \end{aligned}$$

$$c(\mathbf{X}, \mathbf{Z}) = \mathbf{S}_Z \mathbf{V}$$



Quality of representation of variables in PCs

- Check the quality of representation of variable \mathbf{x}^j on PC k

$$I_T(\mathbf{X}) = \sum_{j=1}^p \sum_{s=1}^r r^2(\mathbf{x}^j, \mathbf{z}_s) \quad \text{for normed PCA}$$

- Correlation circle:

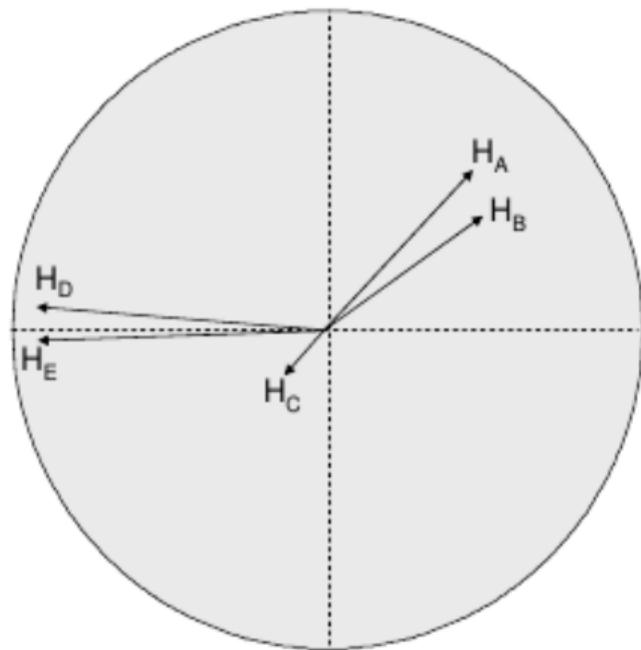
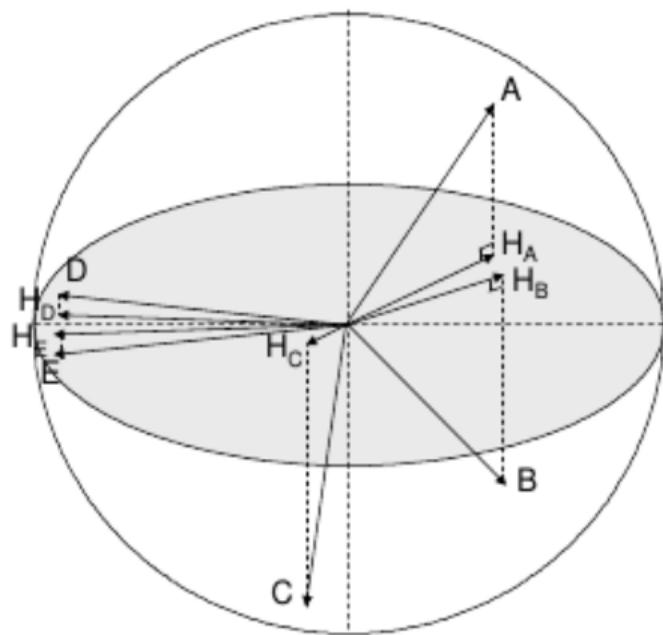
$$\cos^2(\theta\{\mathbf{x}^j, \mathbf{z}_k\}) = \frac{r^2(\mathbf{x}^j, \mathbf{z}_k)}{\sum_{s=1}^r r^2(\mathbf{x}^j, \mathbf{z}_s)}$$

- Only variables with high \cos^2 can be interpreted !
- Contribution of variable \mathbf{x}^j

$$\text{contr}(\mathbf{x}^j, \mathbf{z}_k) = \frac{r^2(\mathbf{x}^j, \mathbf{z}_k)}{\lambda_k}$$

Quality of representation of variables in PCs

Check the quality of representation of variables, close variables are not necessarily similar



Outline

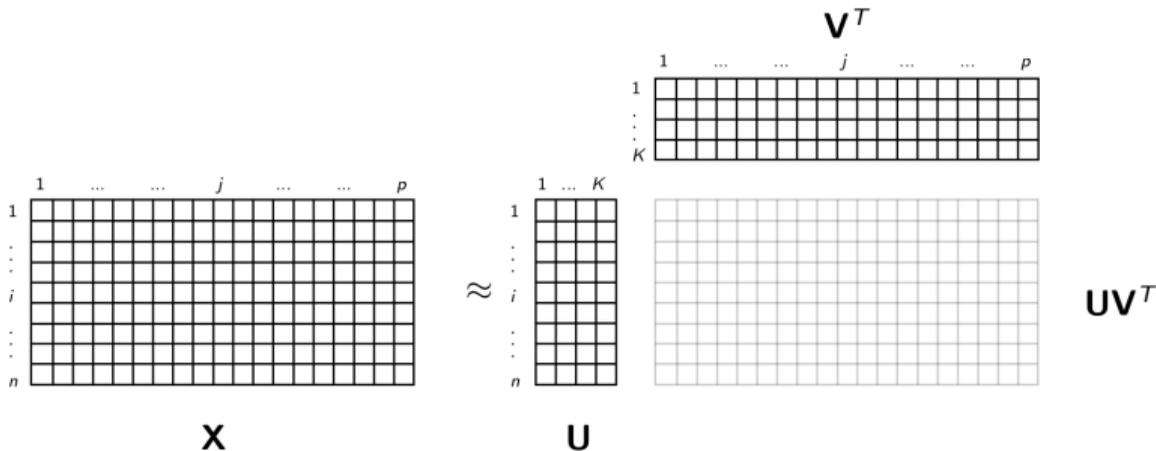
1. Introduction
2. Vectors and distances
3. Defining a new representation
4. Changing Coordinates
5. Dimension Reduction by compression
- 6. Conclusion, extensions**
7. Alternatives to PCA, non linear embedding methods
8. Annexes
9. Principal Components and orthogonal subspaces

Summary

- PCA is the most widely used linear dimension reduction method
- It is based on a change in coordinates to represent the data in a way that preserves the variability of the data
- The new coordinates are provided by the eigenvectors of the empirical variance matrix
- Check the percentage of explained variance and choose the number of components accordingly
- Check the quality of representation of variables to interpret the axes
- Interpret the projection of individuals at the end
- Why does PCA make cluster more visible ?

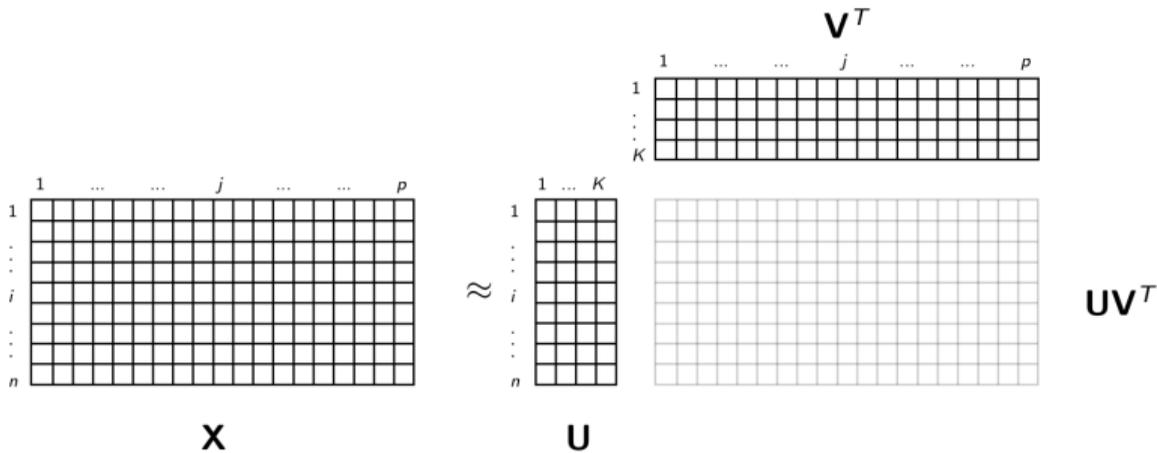
Matrix factorization: $\mathbf{X} \approx \mathbf{UV}^T$

Cells: $\mathbf{U} \in \mathbb{R}^{n \times K}$ }
Genes: $\mathbf{V} \in \mathbb{R}^{p \times K}$ } Low dimensional representation



→ Low-rank representation of \mathbf{X}

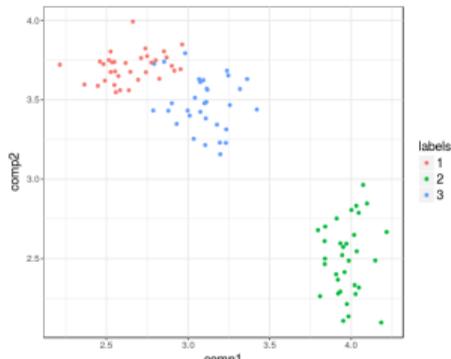
Matrix factorization: $X \approx UV^T$



Data visualization: U

scatter plot $(u_{i1}, u_{i2})_{i=1:n}$

Embeddings

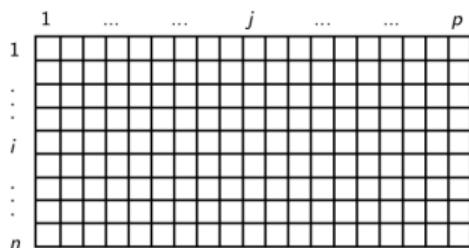


How to interpret the axes ?

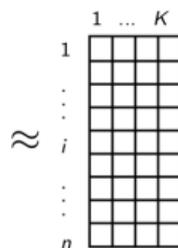
- When genes contribute poorly to axis, their contribution can be put to zero

■ = selected genes ($v_{jk} \neq 0$)

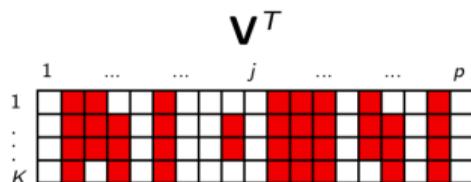
□ = irrelevant genes ($v_{jk} = 0$)



X



U



\approx

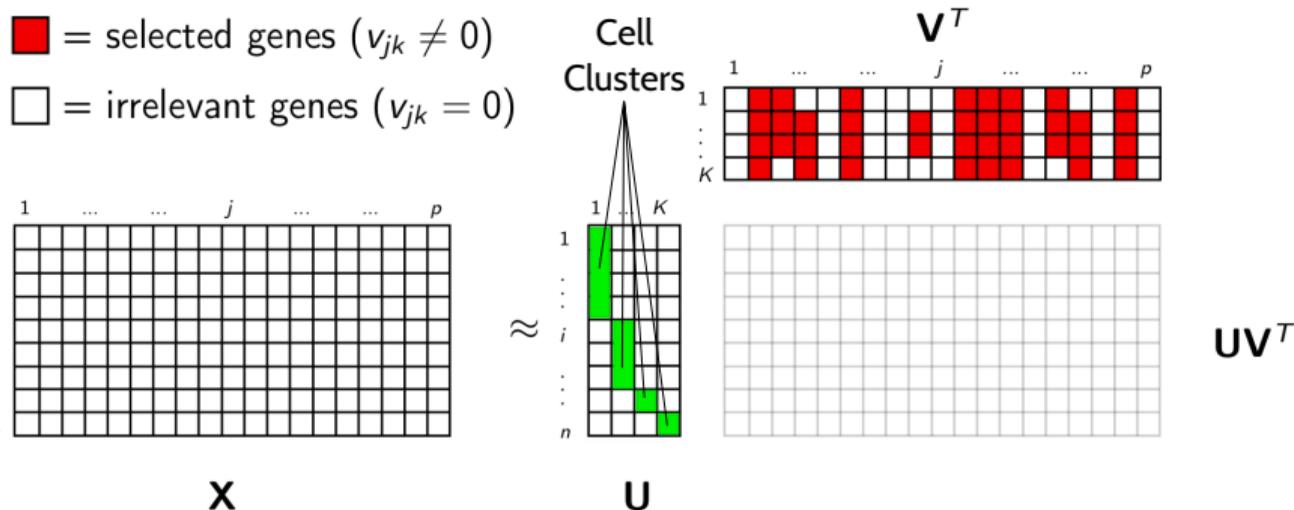


UV^T

- Selected genes can be interpreted in terms of signature.

How to cluster cells in terms of selected variables

- When signatures are selected in \mathbf{V} , this can be used to create clusters of cells in \mathbf{U}



- Compression allows to exhibit variables that make clusters more detectable

Towards embedding methods

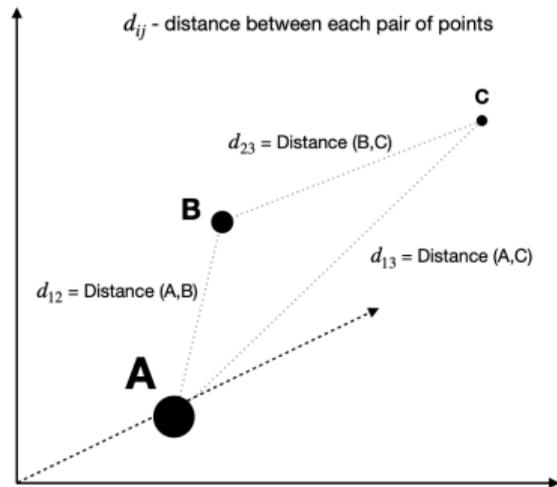
- PCA is based on the duality between the between-variables distance $\mathbf{S} = \mathbf{X}'\mathbf{X}/n$ and the between individuals distance $\mathbf{G} = \mathbf{X}\mathbf{X}'/p$
- \mathbf{U} provides the new coordinates for the individuals
- \mathbf{V} provides the new coordinates for the variables
- Creating a new representation thanks to a linear transform $\mathbf{Z} = \mathbf{X}\mathbf{V}'$ ensures the same transform for each point
- The linear nature of the transform ensures interpretability of PCA
- In the end, data visualization focuses on the representations of individuals, called embeddings.
- Considering embedding allows to extend the notion of dimension reduction to other frameworks

A primer with Multidimensional Scaling

- In many situations only the distance $d_{ii'}$ between individuals (i, i') is available
- The objective of MDS is to find new coordinates $\mathbf{u}_1, \dots, \mathbf{u}_n$ that minimize:

$$\sum_{i, i'} \left(d_{ii'} - \|\mathbf{u}_i - \mathbf{u}_{i'}\|^2 \right)^2$$

- The information regarding the variables is not considered (not available)

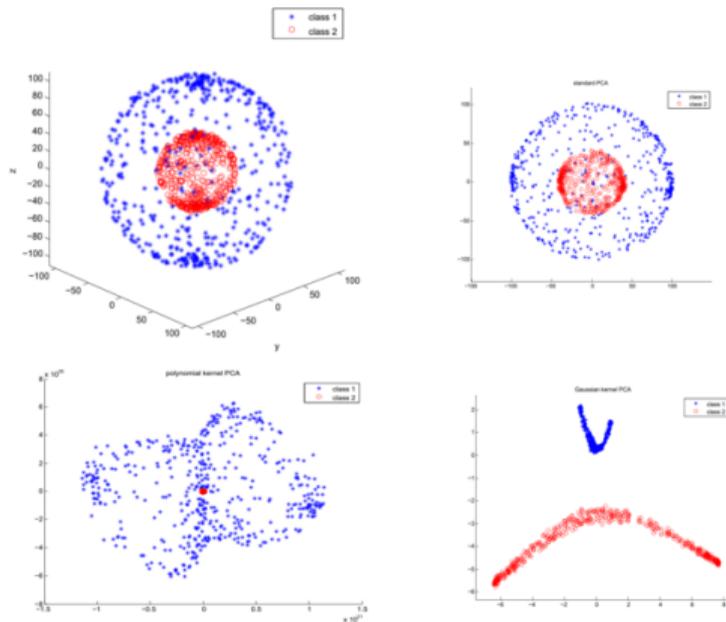


Extending the notion of distance with kernels

- Linear methods are mainly based on euclidean distances
- These distances depend on a dot product
- This dot product can be generalized by the so-called kernel

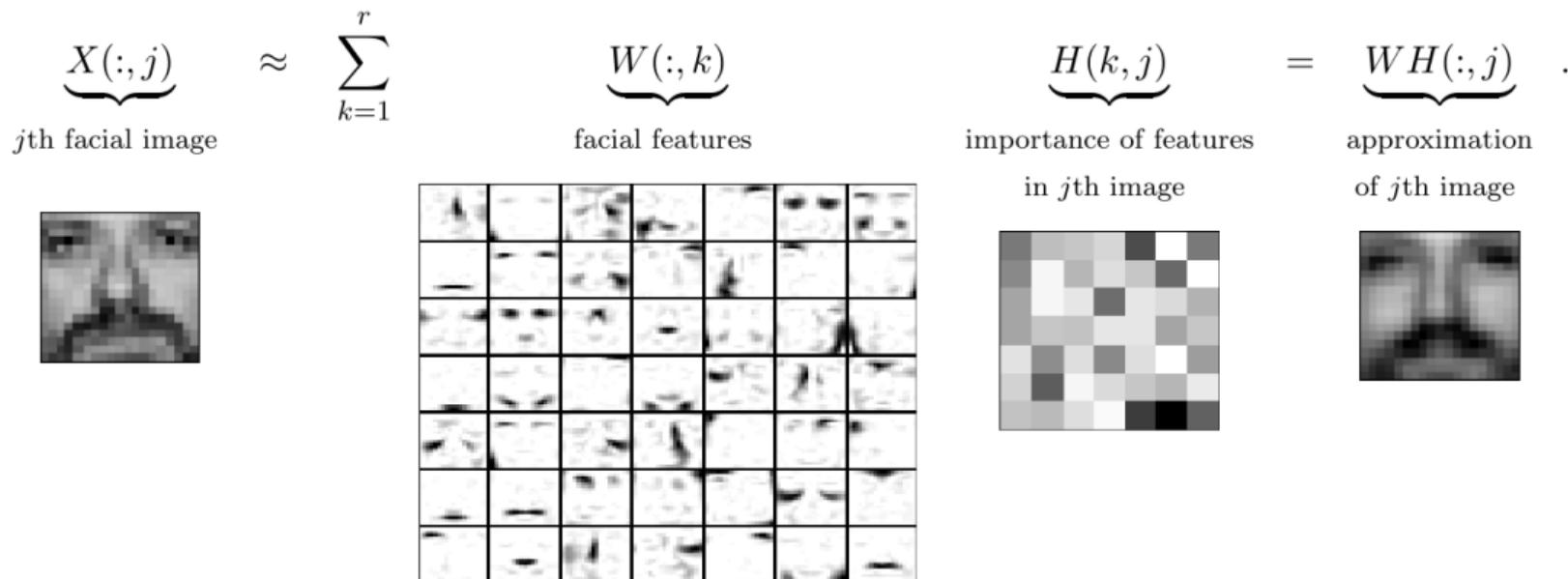
$$K(\mathbf{x}_i, \mathbf{x}_{i'}) = \langle \phi(\mathbf{x}_i), \phi(\mathbf{x}_{i'}) \rangle$$

- ϕ is called the feature map and is unknown
- Grounds most non linear methods (kernel-PCA, kernel MDS, etc)



Accounting for particular characteristics of data

When data are counts, introduce a non-negativity constraint and use NMF

$$\underbrace{X(:, j)}_{j\text{th facial image}} \approx \sum_{k=1}^r \underbrace{W(:, k)}_{\text{facial features}} \underbrace{H(k, j)}_{\text{importance of features in } j\text{th image}} = \underbrace{WH(:, j)}_{\text{approximation of } j\text{th image}} .$$


The diagram illustrates the Non-negative Matrix Factorization (NMF) process for facial image analysis. It shows the decomposition of a facial image $X(:, j)$ into a matrix of facial features $W(:, k)$ and a matrix of feature importance $H(k, j)$. The product $WH(:, j)$ is shown as an approximation of the original image.

The j th facial image is shown as a grayscale portrait of a man's face.

The facial features are shown as a grid of 6x6 small grayscale images, each representing a different feature (e.g., eyes, nose, mouth, etc.).

The importance of features in the j th image is shown as a 6x6 grayscale heatmap, where darker pixels indicate higher importance.

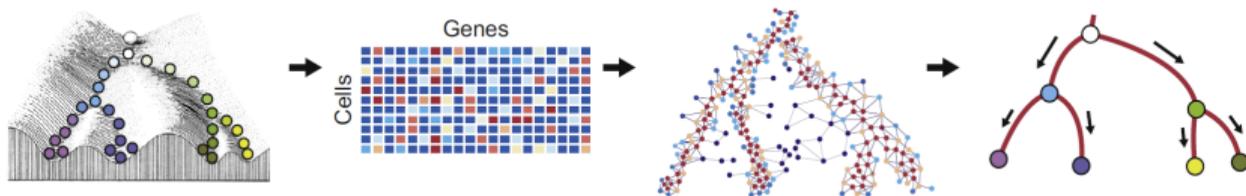
The approximation of the j th image is shown as a grayscale portrait of the same man's face, reconstructed from the features and their importance weights.

Outline

1. Introduction
2. Vectors and distances
3. Defining a new representation
4. Changing Coordinates
5. Dimension Reduction by compression
6. Conclusion, extensions
- 7. Alternatives to PCA, non linear embedding methods**
8. Annexes
9. Principal Components and orthogonal subspaces

Beyond Linear projections

- Linear methods are powerful for planar structures
- High dimensional datasets are characterized by multiscale properties (local / global structures)
- May not be the most powerful for manifolds
- Non Linear projection methods aim at preserving local characteristics of distances



Stochastic Neighbor Embedding [?]

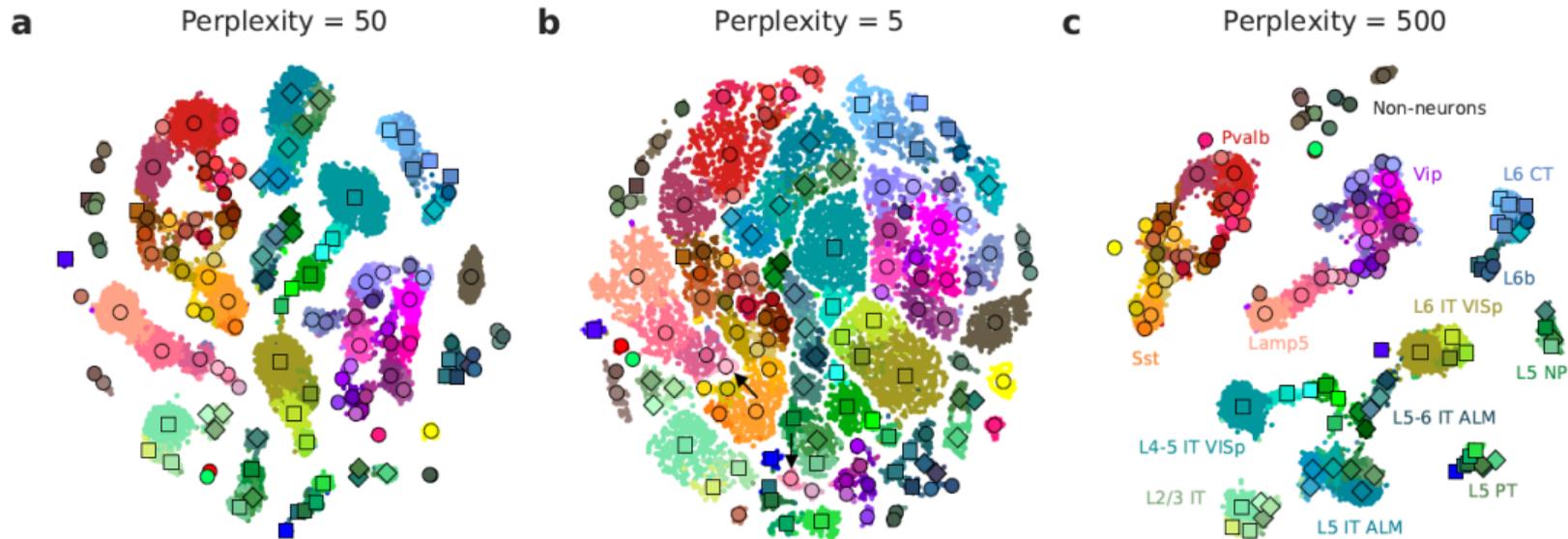
- (x_1, \dots, x_n) are the points in the high dimensional space \mathbb{R}^p ,
- Consider a similarity between points:

$$p_{i|j} = \frac{\exp(-\|x_i - x_j\|^2/2\sigma_i^2)}{\sum_{k \neq i} \exp(-\|x_k - x_j\|^2/2\sigma_k^2)}, \quad p_{ij} = (p_{i|j} + p_{j|i})/2N$$

- σ smooths the data (linked to the regularity of the target manifold)
- σ is chosen such that the entropy of p is fixed to a given value of the so-called perplexity

$$\exp\left(-\sum_{ij} p_{ij} \log(p_{ij})\right)$$

Visual inspection of the influence of σ [?]



tSNE and Student / Cauchy kernels

- Consider (y_1, \dots, y_n) are points in the low dimensional space \mathbb{R}^2
- Consider a similarity between points in the new representation:

$$q_{ij} = \frac{\exp(-\|y_i - y_j\|^2)}{\sum_{k \neq i} \exp(-\|y_k - y_j\|^2)}$$

- Robustify this kernel by using Student(1) kernels (ie Cauchy)

$$q_{ij} = \frac{(1 + \|y_i - y_j\|^2)^{-1}}{\sum_{k \neq i} (1 + \|y_i - y_k\|^2)^{-1}}$$

Optimizing tSNE

- Minimize the KL between p and q so that the data representation minimizes:

$$C(y) = \sum_{ij} KL(p_{ij}, q_{ij})$$

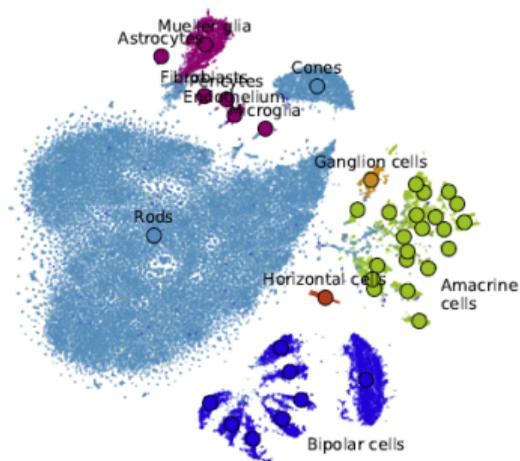
- The cost function is not convex

$$\left[\frac{\partial C(y)}{\partial y} \right]_i = \sum_j (p_{ij} - q_{ij})(y_i - y_j)$$

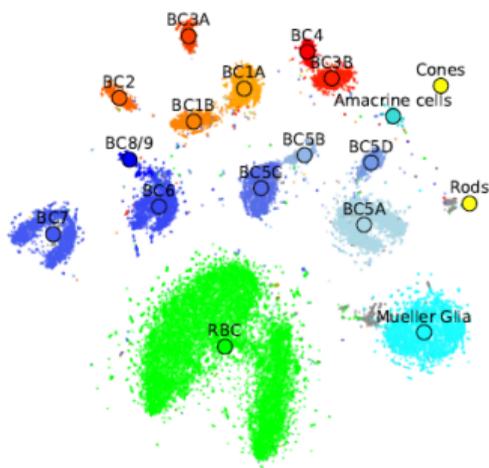
- Interpreted as the resultant force created by a set of springs between the map point y_i and all other map points $(y_j)_j$. All springs exert a force along the direction $(y_i - y_j)$.
- $(p_{ij} - q_{ij})$ is viewed as a stiffness of the force exerted by the spring between y_i and y_j .

tSNE examples on single cell RNASeq data 1 [?]

a Macosko et al. 2015



b Shekhar et al. 2016



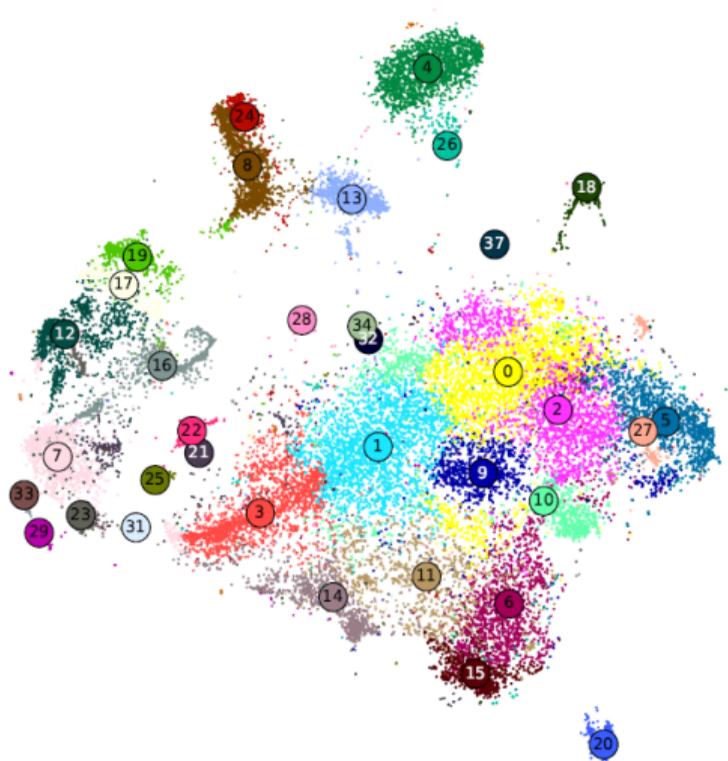
c Harris et al. 2018



tSNE examples on single cell RNASeq data 1 [?]

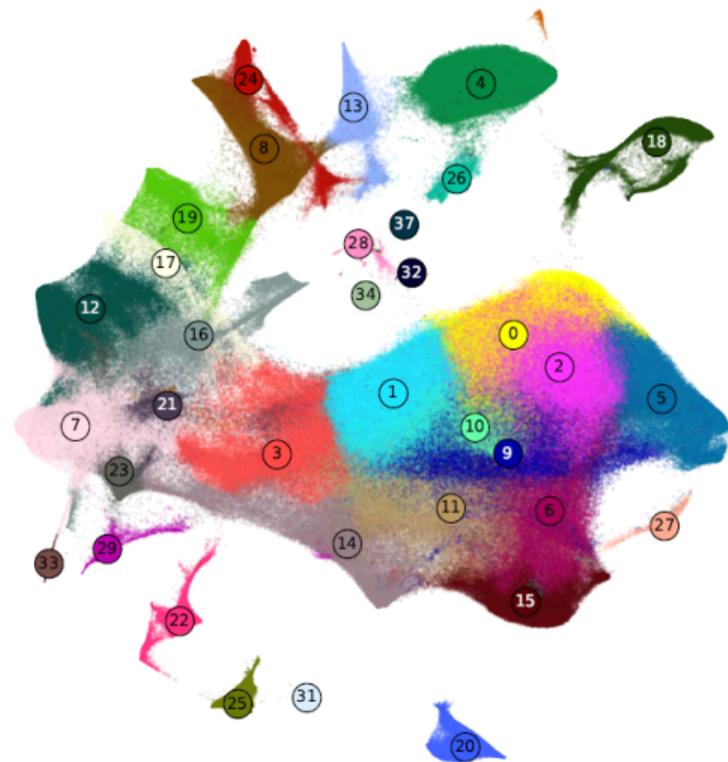
a

$N = 25\,000$

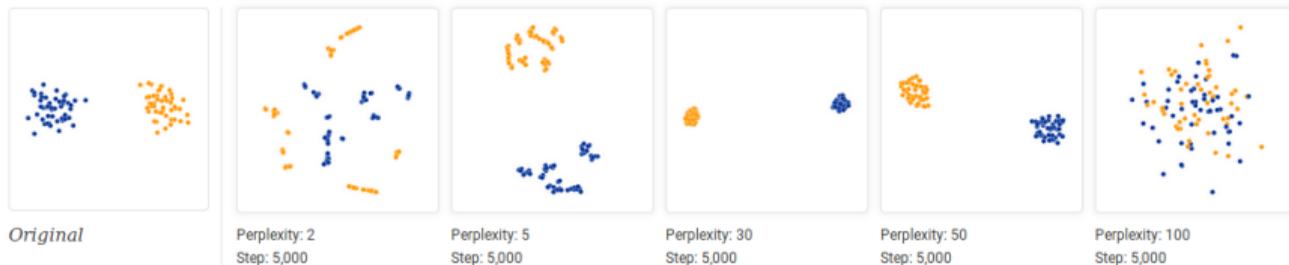


b

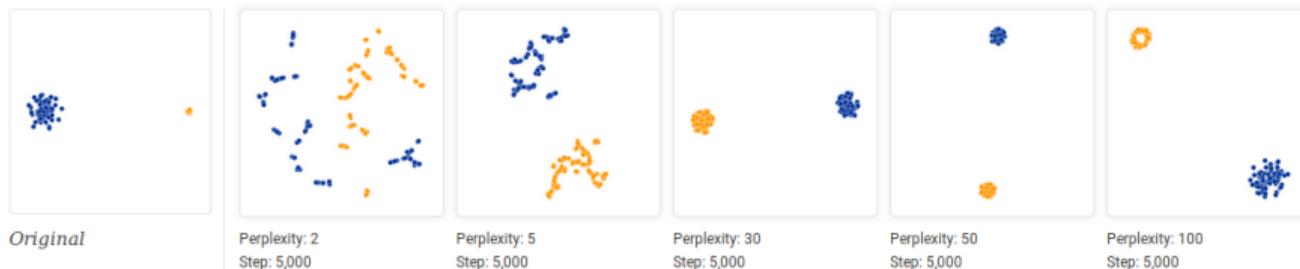
$N = 1\,306\,127$



Effect of Hyperparameters : Perplexity

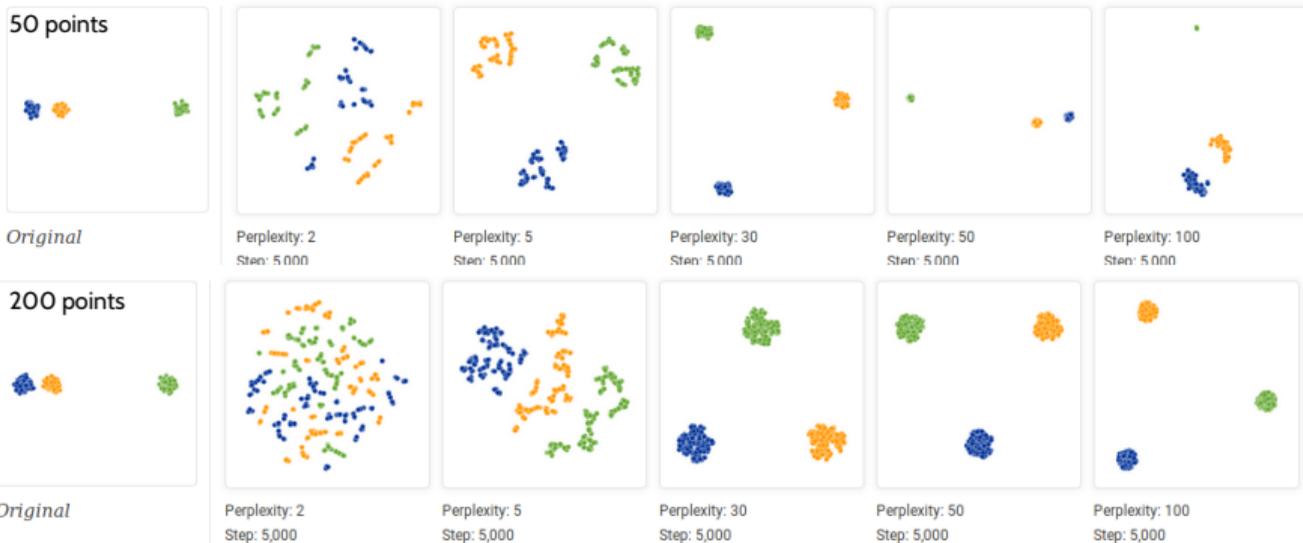


tSNE does not account for heteroscedasticity

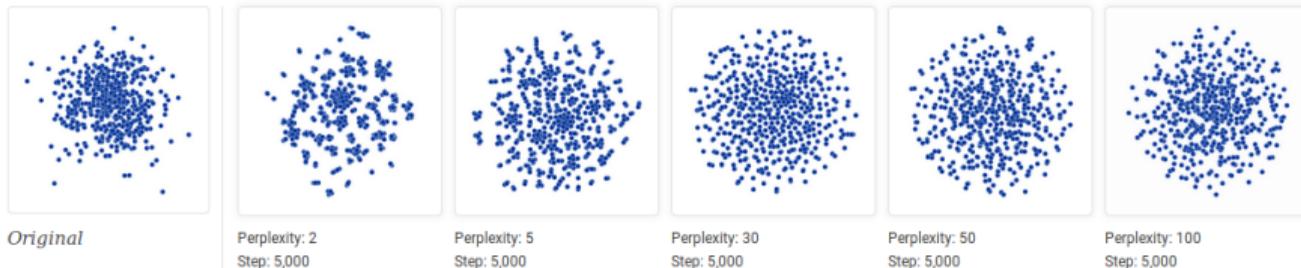


<https://distill.pub/2016/misread-tsne/>

tSNE does not account for between-cluster distance



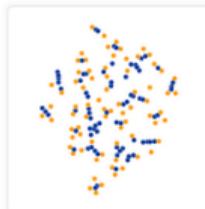
What about random noise ?



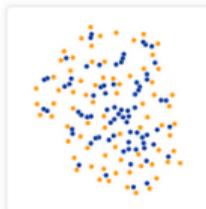
Catching Complex Geometries



Original



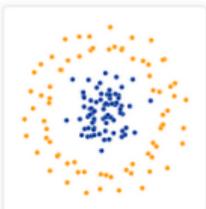
Perplexity: 2
Step: 5,000



Perplexity: 5
Step: 5,000



Perplexity: 30
Step: 5,000



Perplexity: 50
Step: 5,000



Perplexity: 100
Step: 5,000



Original



Perplexity: 2
Step: 5,000



Perplexity: 5
Step: 5,000



Perplexity: 30
Step: 5,000



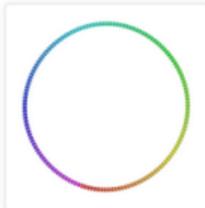
Perplexity: 50
Step: 5,000



Perplexity: 100
Step: 5,000



Original



Perplexity: 2
Step: 5,000



Perplexity: 5
Step: 5,000



Perplexity: 30
Step: 5,000



Perplexity: 50
Step: 5,000

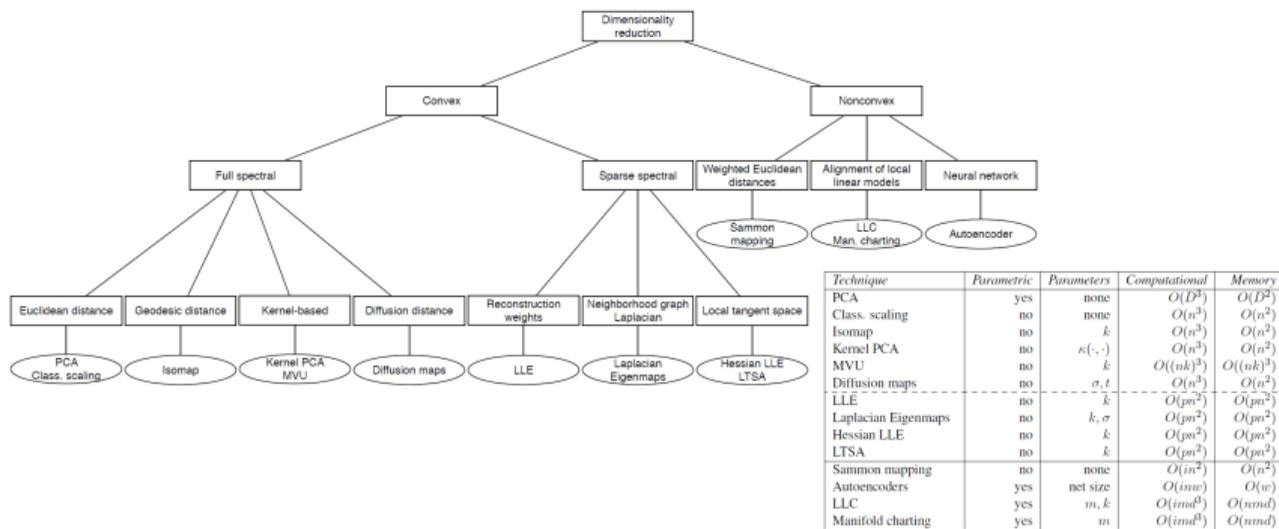


Perplexity: 100
Step: 5,000

Properties of t-SNE

- good at preserving local distances (intra-cluster variance)
- not so good for global representation (inter-cluster variance)
- hence good at creating clusters of points that are close, but bad at positioning clusters wrt each other
- preprocessing very important : initialize with PCA and feature selection plus log transform (non linear transform)
- percent of explained variance ? interpretation of the q distribution ?

A taxonomy of Dimension Reduction Methods [?]



Conclusions of a comparative study [?]

- local methods suffer from the choice of the smoothing (neighborhood) parameter
- Kernel PCA suffers from the choice of the Kernel to correctly approximate the manifold.
- Setting the optimization problem is the key (convex or not), trivial solutions, local optima, computationally feasible
- nonlinear techniques for dimensionality reduction are, despite their large variance, often not capable of outperforming traditional linear techniques such as PCA.

Useful links

- <https://towardsdatascience.com/>
- PCA for datascience
- Link to a tuto on dot products
- Wiki for Linear Transforms
- Book for the introduction to machine learning (C.-A. Azencott)
- Book for the introduction to machine learning (James, Witten, Hastie, Tibshirani)
- PCA in ecology <http://pbil.univ-lyon1.fr/ade4/>
- PCA in general http://factominer.free.fr/index_fr.html

References

- [1] LJP Van der Maaten, EO Postma, and HJ Van den Herik. Dimensionality reduction: A comparative review. *TiCC*, 2009.
- [2] J. G. Greener, S. M. Kandathil, L. Moffat, and D. T. Jones. A guide to machine learning for biologists. *Nat Rev Mol Cell Biol*, 23(1):40–55, 01 2022.
- [3] Dmitry Kobak and Philipp Berens. The art of using t-sne for single-cell transcriptomics. *bioRxiv*, 2018.
- [4] M. Ringnér. What is principal component analysis? *Nat Biotechnol*, 26(3):303–304, Mar 2008.
- [5] Laurens van der Maaten and Geoffrey Hinton. Visualizing Data using t-SNE. *Journal of Machine Learning Research*, 9(Nov):2579–2605, 2008.

Outline

1. Introduction
2. Vectors and distances
3. Defining a new representation
4. Changing Coordinates
5. Dimension Reduction by compression
6. Conclusion, extensions
7. Alternatives to PCA, non linear embedding methods
8. Annexes
- 9. Principal Components and orthogonal subspaces**

Expectation / Variance for matrices

- Given $Y_i \in \mathbb{R}^p$, $A \in \mathbb{R}^{q \times p}$,

$$\mathbb{E}(AY_i) = A \times \mathbb{E}(Y_i),$$

- The variance of a linear combination of Y

$$\mathbb{V}(AY_i) = A\mathbb{V}_{p \times p}(Y_i)A',$$

Outline

1. Introduction
2. Vectors and distances
3. Defining a new representation
4. Changing Coordinates
5. Dimension Reduction by compression
6. Conclusion, extensions
7. Alternatives to PCA, non linear embedding methods
8. Annexes
9. Principal Components and orthogonal subspaces

Decomposition of \mathbb{R}^p into orthogonal subspaces

- Let us consider p orthogonal subspaces $(E_k)_{k=1,p}$ each subspace spanned by an individual axis (dim 1):

$$\mathbb{R}^p = \bigoplus_{k=1}^p E_k,$$

- Orthogonal projection of $X_i \in \mathbb{R}^p$ on a subspace $E_k = \text{vect}(Z_k)$

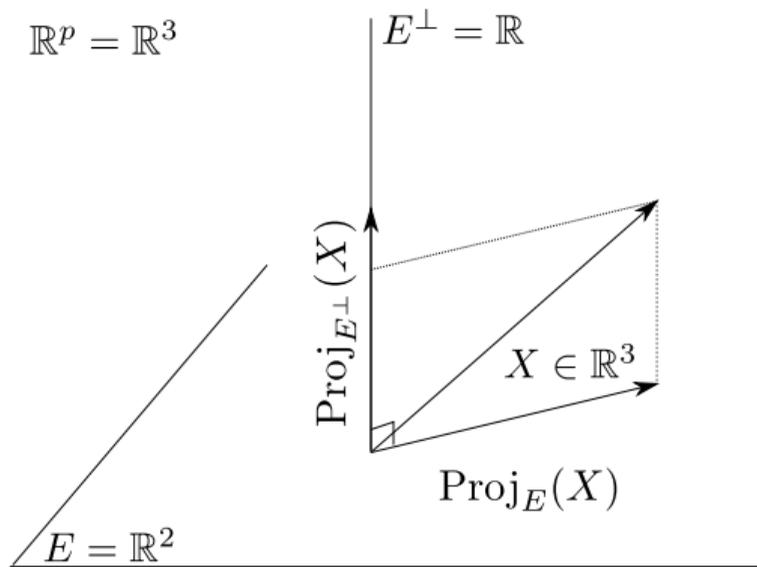
$$\text{Proj}_{E_k}(X_i) = X_i V_k \in \mathbb{R}$$

- The inertia of X wrt E_k measures the proximity of E_k from X

$$I_{E_k}(X) = \frac{1}{n} \sum_{i=1}^n \|X_i - \text{Proj}_{E_k}(X_i)\|_2^2$$

- Let E_k^\perp denotes the orthogonal complement of subspace E_k .

Pythagore - Huyguens Theorem



$$I_T(X) = I_E(X) + I_{E^\perp}(X) = I(\text{Proj}_E(X)) + I(\text{Proj}_{E^\perp}(X))$$

Construction of principal components (PC)

- Resume the data X by a new dataset $Z_{n \times K}$, $K \leq p$ and K fixed
- The new axis spans the 1-dim subspaces $\left(E_k = \text{vect}(Z_k)\right)_k$

$$\forall k, k', \quad E_k \perp E_{k'}$$

- $Z = [Z_1, \dots, Z_K]$ constitute independent PCs (easy interpretation)
- $Z_k \in \mathbb{R}^n$ is defined as a linear combination of the variables

$$Z_k = XV_k, \quad V_k = (V_{jk})_j \in \mathbb{R}^p$$

- $V_{p \times K} = [V_1, \dots, V_K]$ is the matrix of contributions (weights) of variables $(X^j)_j$

$$Z_{n \times K} = X_{n \times p} V_{p \times K}$$

Decomposition of the Inertia on the PCs

$$\begin{aligned}I_T(X) &= \frac{1}{n} \sum_{i=1}^n \sum_{k=1}^p \|X_i - \text{Proj}_{E_k}(X_i) + \text{Proj}_{E_k}(X_i)\|^2 \\&= \frac{1}{n} \sum_{i=1}^n \sum_{k=1}^p \|X_i - \text{Proj}_{E_k}(X_i)\|^2 + \frac{1}{n} \sum_{i=1}^n \sum_{k=1}^p \|\text{Proj}_{E_k}(X_i)\|^2 \\&= \frac{1}{n} \sum_{i=1}^n \sum_{k=1}^p \|X_i - Z_{ik}\|^2 + \frac{1}{n} \sum_{i=1}^n \sum_{k=1}^p \|Z_{ik}\|^2 \\&= \frac{1}{n} \sum_{i=1}^n \sum_{k=1}^p \|X_i - X_i V_k\|^2 + \frac{1}{n} \sum_{i=1}^n \sum_{k=1}^p \|X_i V_k\|^2\end{aligned}$$

Orthogonal Components with maximal variance

- We want to resume the variability of the dataset
- Find the PCs that explain the maximum of the observed variance:

$$\frac{1}{n} \sum_{i=1}^n \|\text{Proj}_{E_k}(X_i)\|^2 = \frac{1}{n} \sum_{i=1}^n \|Z_{ik}\|^2 = \frac{1}{n} V_k' (X'X) V_k = \frac{1}{n} V_k' \Sigma V_k$$

- The optimization scheme is iterative, and for the k th PC:

$$\hat{V}_k = \arg \max_{V \in \mathbb{R}^p, \|V\|_2=1} \left(\frac{1}{n} V' X' X V \right) \quad \text{with } Z_k \perp (Z_1, \dots, Z_{k-1})$$

Constrained optimization

- To account for the orthogonality constraint, we introduce the Lagrange multipliers

$$\mathcal{L}(V, \mu) = \frac{1}{n} V' X' X V - \mu (V' V - 1)$$

$$\frac{\partial \mathcal{L}}{\partial \mu} = V' V - 1$$

$$\frac{\partial \mathcal{L}}{\partial V} = 2X' X V - \mu V$$

- Which gives the following solution

$$\begin{aligned} V' V &= 1 \\ X' X V &= \mu V \end{aligned}$$

- The optimal solution is provided by the eigenvectors of the covariance matrix Σ

Spectral decomposition of symmetric real matrices

- Let $A \in \mathbb{R}^{n,n}$ a symmetric real matrix
- Spectral decomposition theorem: there exists $\lambda_1 \geq \dots \geq \lambda_n \in \mathbb{R}$ and an orthogonal basis $\{U_1, \dots, U_n\}$ of \mathbb{R}^n such that

$$A = \sum_{k=1}^n \lambda_k U_k U_k'$$

- The spectral decomposition can also be written:

$$A = U \operatorname{diag}(\lambda_1, \dots, \lambda_n) U'$$

Positive Semi-Definite Matrices

- A symmetric real matrix is positive semi-definite (sdp) if

$$\forall x \in \mathbb{R}^n, x'Ax \geq 0$$

- Semi-Definite positiveness is equivalent to $\lambda_1 \geq \dots \geq \lambda_n \geq 0$, since

$$x'Ax = \sum_{k=1}^n \lambda_k \langle x, U_k \rangle^2$$

- For any $n \times p$ matrix A , the matrices $A'A$ and AA' are symmetric positive semidefinite

Singular Value Decomposition Theorem

- Any matrix $A \in \mathbb{R}^{n,p}$ of rank r can be decomposed as

$$A = \sum_{k=1}^r \mu_k U_k V_k'$$

- $r = \text{rank}(A)$
- $\mu_1 \geq \dots \geq \mu_r > 0$
- $\{\mu_1^2, \dots, \mu_r^2\}$ are the non-zero eigenvalues of $A'A$ and of AA'
- $\{\mu_1, \dots, \mu_r\}$ are called the singular values of A
- $\{U_1, \dots, U_r\}$ and $\{V_1, \dots, V_r\}$ are two orthonormal families of \mathbb{R}^n and \mathbb{R}^p such that:

$$AA'U_k = \mu_k^2 U_k, \quad A'AV_k = \mu_k^2 V_k$$

Singular Value Decomposition of $X'X$ and XX'

- (U_1, \dots, U_K) , the eigen vectors of the Gram matrix

$$G_{n \times n} = \frac{1}{p} XX' = \frac{1}{p} \sum_{k=1}^K \mu_k^2 U_k U_k'$$

- (V_1, \dots, V_K) , the eigen vectors of the covariance matrix

$$\Sigma_{p \times p} = \frac{1}{n} X'X = \frac{1}{n} \sum_{k=1}^K \mu_k^2 V_k V_k'$$

- Then we have

$$\begin{aligned}(XX')U_k &= \mu_k X V_k = \mu_k^2 U_k \\(X'X)V_k &= \mu_k X' U_k = \mu_k^2 V_k\end{aligned}$$

Low-rank approximation of X (1)

- $X \in \mathbb{R}^{n,p}$, s.t. $\text{rank}(X) = r$, there exists
 - $\mu_1 \geq \dots \geq \mu_r > 0$, with $D = \text{diag}(\mu_1, \dots, \mu_r)$,
 - $\{\mu_1, \dots, \mu_r\}$, are the singular values of X
 - two orthogonal matrices $\tilde{U} \in \mathbb{R}^{n \times r}$ and $\tilde{V} \in \mathbb{R}^{p \times r}$ with

$$\tilde{U}'\tilde{U} = I_r, \quad \tilde{V}'\tilde{V} = I_r,$$

$$U = \tilde{U}D, \quad V = \tilde{V}D,$$

- Such that

$$X = UV' = \tilde{U}D\tilde{V}' = \sum_{k=1}^r \mu_k \tilde{U}_k \tilde{V}_k'$$

- Then we have

$$X'\tilde{U}_k = \mu_k \tilde{V}_k, \quad X\tilde{V}_k = \mu_k \tilde{U}_k$$

Low Rank approximation of X (2)

- If $\text{rank}(X) = r$ (unknown), in practice we choose $K \leq p$ to provide a "low-rank" approximation of X .
- Denoting $\hat{X}_K = U_{1:K} V'_{1:K}$ this approximation of $\text{rank}(\hat{X}_K) = K$
- PCA can be restated as the approximation of X st

$$\|X - \hat{X}_K\|_F^2 = \min_{B \in \mathcal{M}_{n,K}, \text{rk}(B)=K} \|X - B\|_F^2 = \sum_{k=K+1}^r \mu_k^2$$

- PCA provides the best low-rank approximation for the Frobenius norm

$$\hat{X}_K = \arg \min_{B \in \mathcal{M}_{n,K}, \text{rk}(B)=K} \|X - B\|_F^2$$